



**Weill Cornell
Medicine**



Weill Cornell Medicine
Institute of Artificial Intelligence
for Digital Health

Disentangling the Clinical Complexity of Heterogeneous Diseases Through Data-Driven Subphenotyping with Large Scale Electronic Health Records

Fei Wang

Associate Professor, Department of Population Health Sciences

Director, Institute of Artificial Intelligence for Digital Health

Weill Cornell Medicine, Cornell University

few2001@med.cornell.edu

 @feiwang03



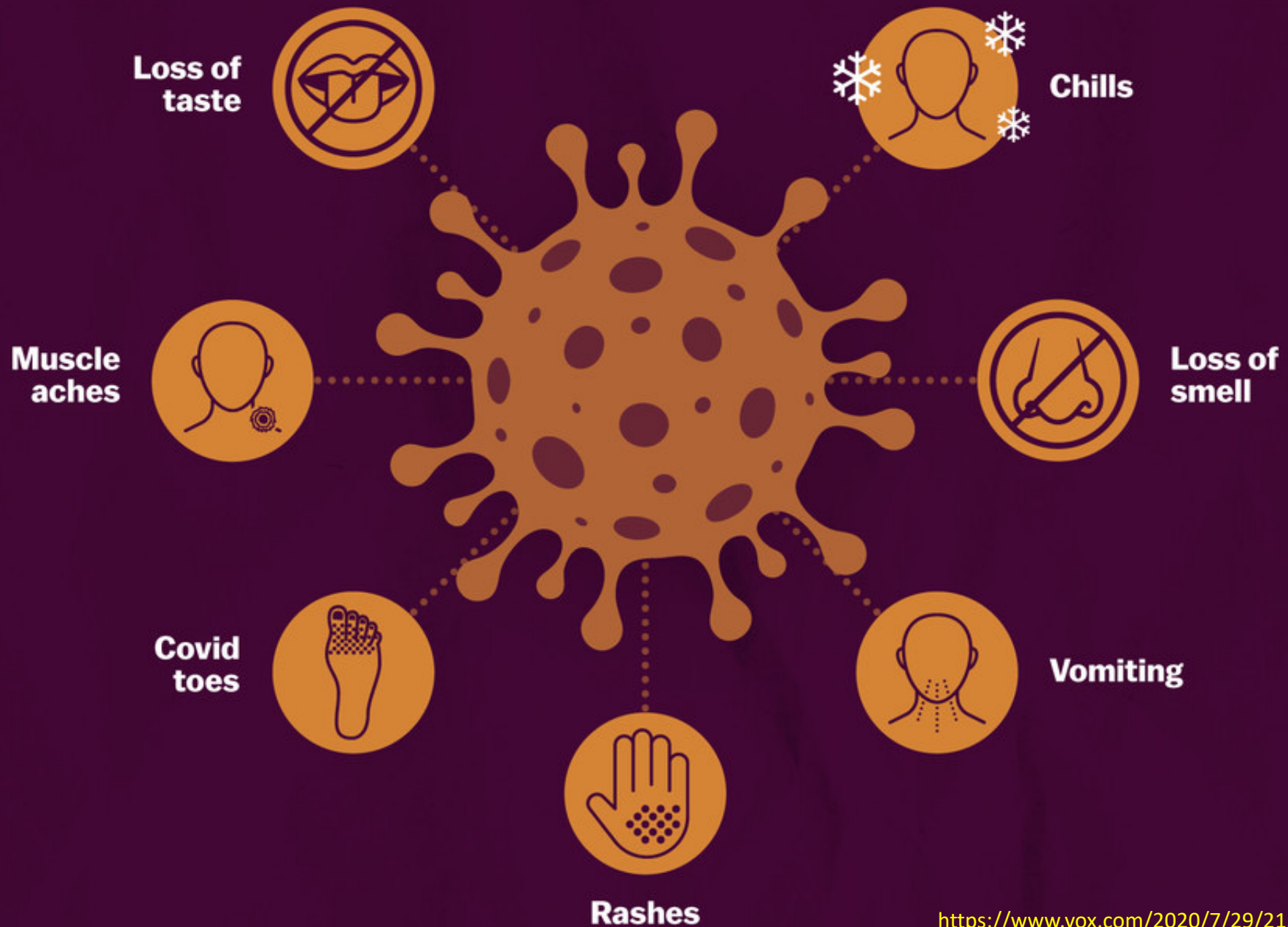
<https://wcm-wanglab.github.io/index.html>

Outline

- Introduction
- Subphenotyping of COVID-19 at infection confirmation
- Subphenotyping of Severe COVID-19 after Mechanical Ventilation
- Subphenotyping of Long COVID
- Discussions

Outline

- Introduction
- Subphenotyping of COVID-19 at infection confirmation
- Subphenotyping of Severe COVID-19 after Mechanical Ventilation
- Subphenotyping of Long COVID
- Discussions

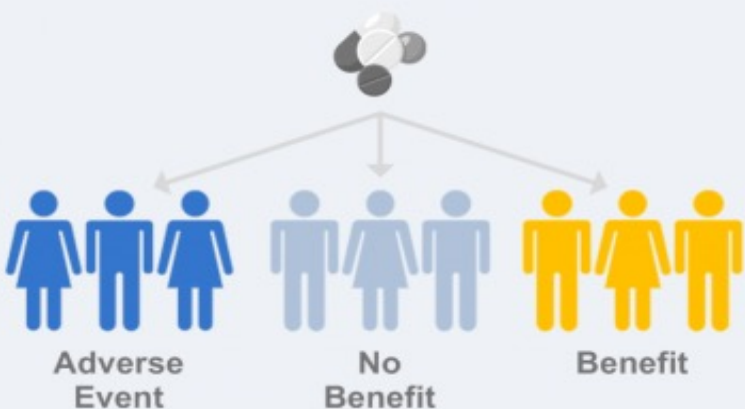




Traditional Medicine



Therapy (mainly Rx)



Stratified Medicine

Patients are grouped by:

- Disease Sub-types
- Risk Profiles
- Demographics
- Socio-economic Factors
- Clinical Features
- Biomarkers
- Molecular Sub-populations



Therapy (mainly Rx)



Precision Medicine

Individual patient level:

- Genomics and Omics
- Lifestyle
- Preferences
- Health History
- Medical Records
- Compliance
- Exogenous Factors



Companion Diagnostic (CDx) Biomarker

Therapy (Rx + Dx = CDx)

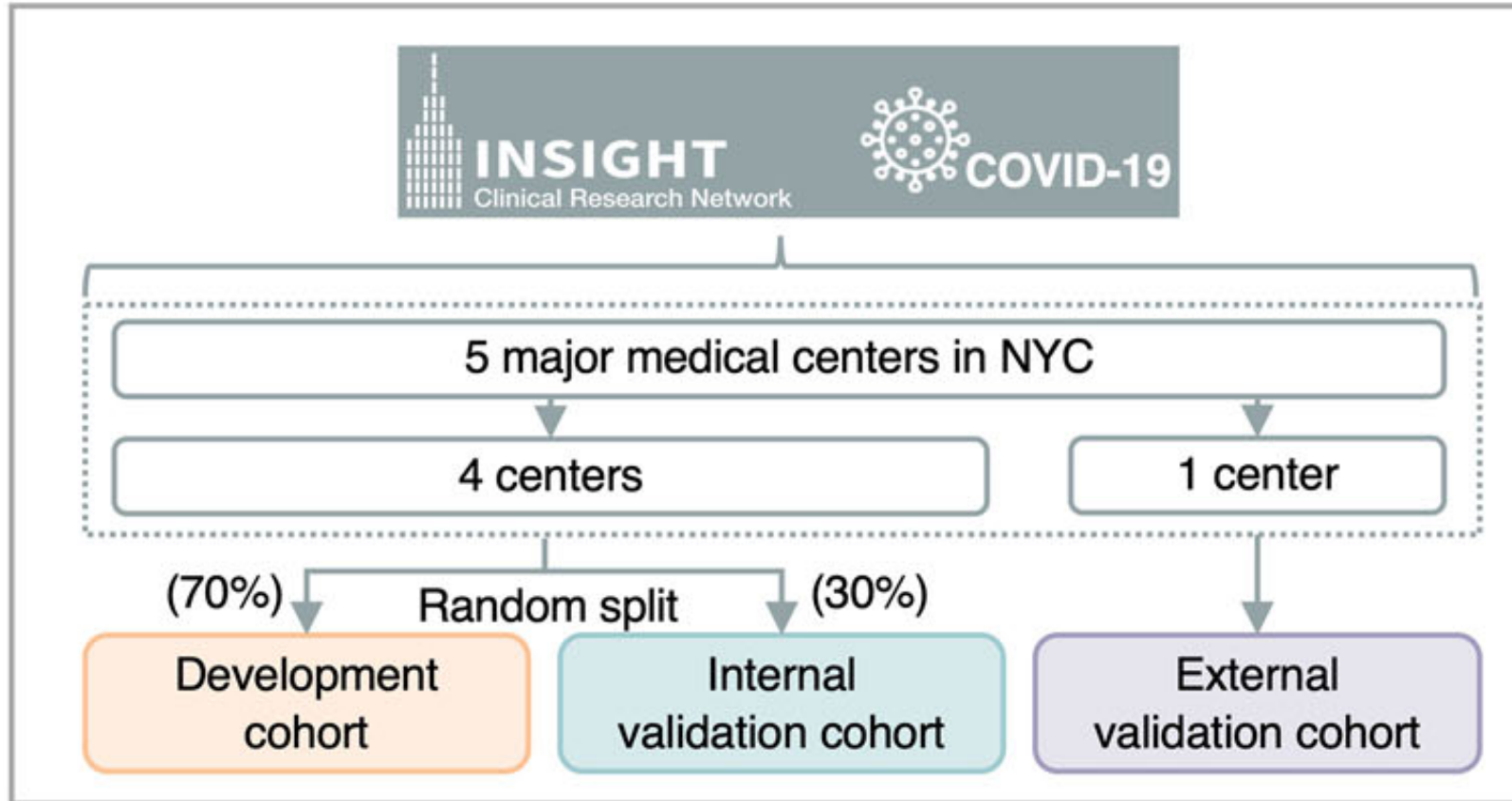


Precision medicine research enables development and delivery of the right patient intervention

Outline

- Introduction
- **Subphenotyping of COVID-19 at infection confirmation**
- Subphenotyping of Severe COVID-19 after Mechanical Ventilation
- Subphenotyping of Long COVID
- Subphenotyping of AD
- Discussions

Overall Setup

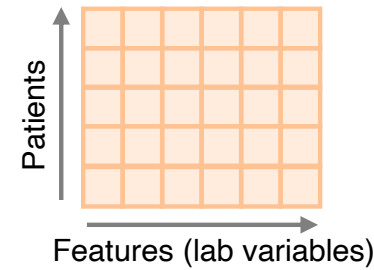


Data preparation

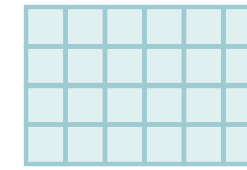


- Collection of presenting laboratory test data;
- Data scaling;
- Imputation of missing data

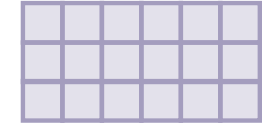
Development cohort



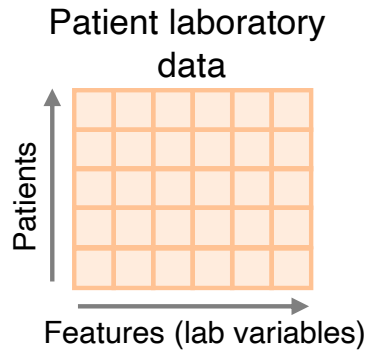
Internal validation cohort



External validation cohort

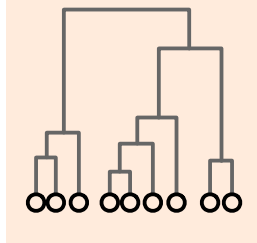


Derivation of subphenotypes



Development cohort

Agglomerative hierarchical clustering



...

Subphenotypes (i.e., patient subgroups)

Sensitivity analysis

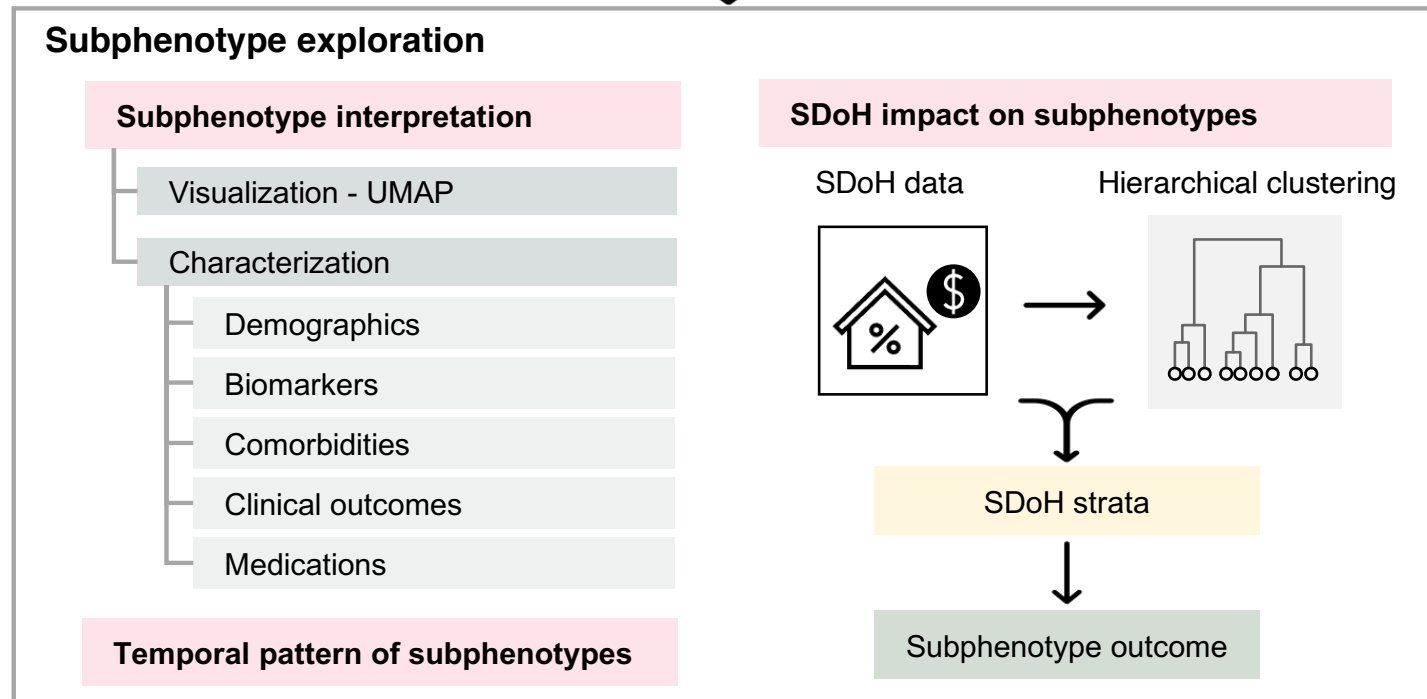
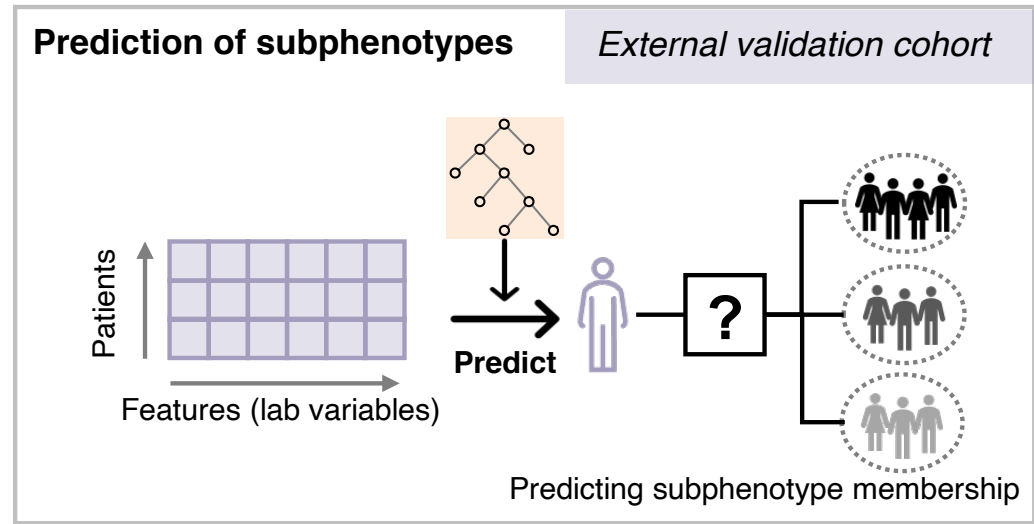
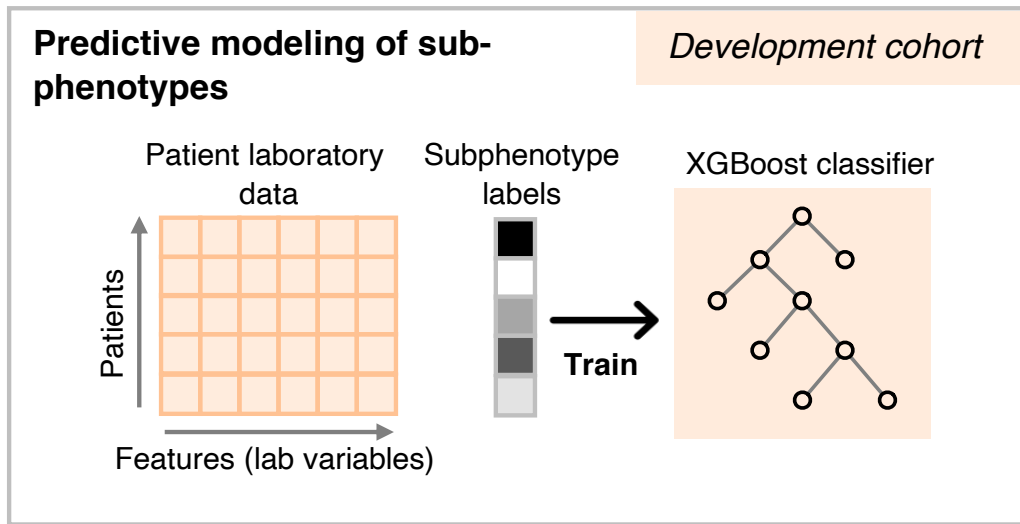
Development cohort

1. Sensitivity to outliers
Re-deriving subphenotypes after removing outliers.
2. Sensitivity to clustering methods
Re-deriving subphenotypes using another clustering method - Gaussian mixture model.

Re-derivation

Internal validation cohort

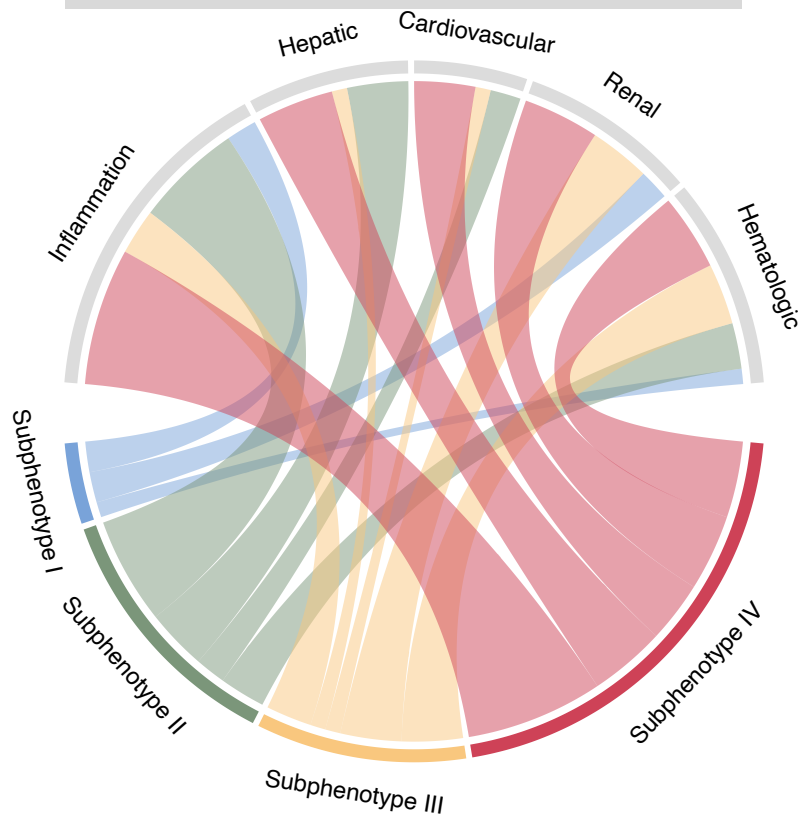
Re-deriving subphenotypes in internal validation cohort using agglomerative hierarchical clustering.



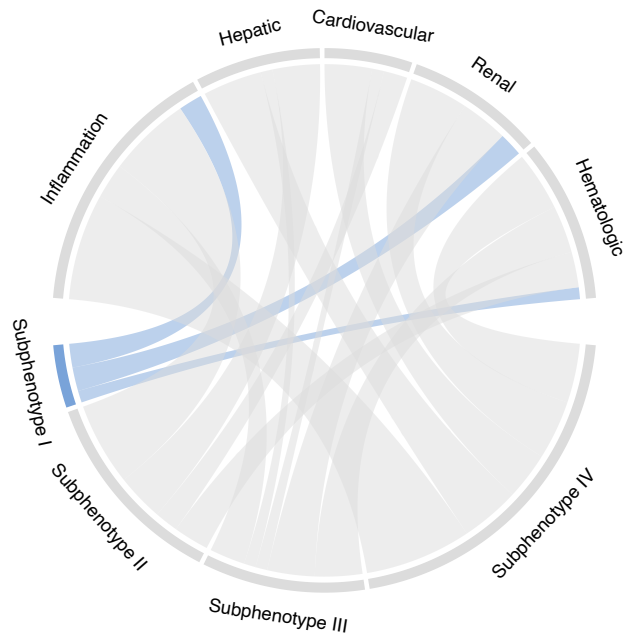
Cohort Characteristics

Characteristics	Cohort		
	Development cohort	Internal validation cohort	External validation cohort
No. of patients	8,199	3,519	2,700
Construction method	70% patients (randomly selected) from 4 sites of INSIGHT network: NYU-LMC, NYP-WCMC, MSHS, and MMC	Remaining 30% patients from 4 sites of INSIGHT network: NYU Langone Medical Center, NYP-WCMC, Mount Sinai Health System, and Montefiore Medical Center	NYP-CUMC
Age, y, Median (IQR)	63.53 [50.57 - 75.15]	63.51 [50.95 - 75.17]	65.58 (51.08 - 77.39)
Sex female, N (%)	3,787 (46.2)	1,585 (45.0)	1,305 (48.3)
Race, N (%)			
White	2,036 (24.8)	838 (23.8)	675 (25.0)
Black	2,155 (26.3)	915 (26.0)	545 (20.2)
Asian	409 (5.0)	193 (5.5)	28 (1.0)
Multiple race	39 (0.5)	13 (0.4)	912 (33.8)
Other/unknown	3560 (43.4)	1560 (44.3)	540 (20.0)
Outcomes (60 days), N (%)			
Mortality	1529 (18.65)	696 (19.78)	556 (20.59)
Mechanical ventilation (intubation)	1154 (14.07)	497 (14.12)	248 (9.19)
ICU admission	1494 (18.22)	661 (18.78)	-

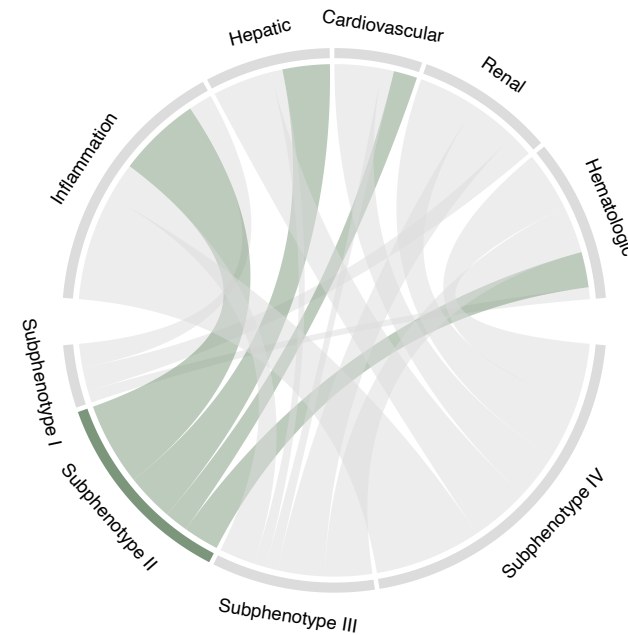
Abnormal biomarkers vs. subphenotypes



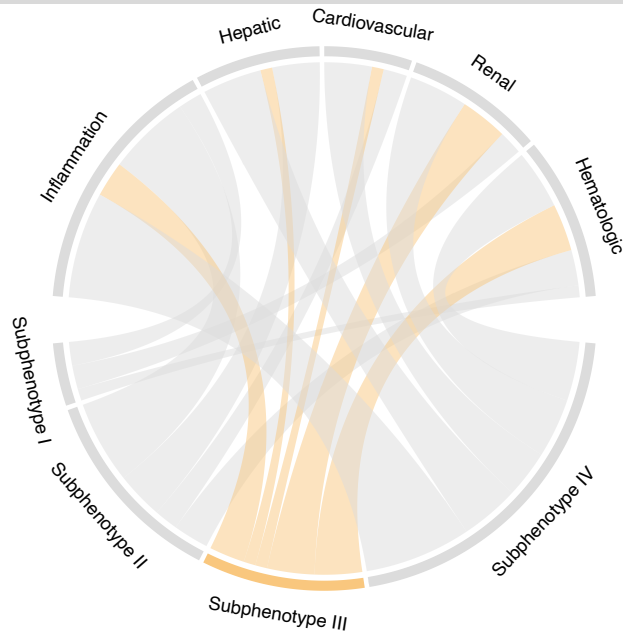
Abnormal biomarkers vs. Subphenotype I



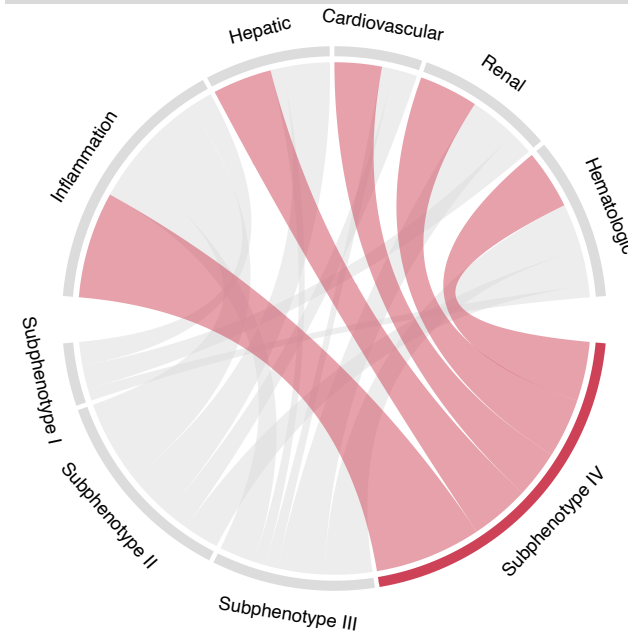
Abnormal biomarkers vs. Subphenotype II



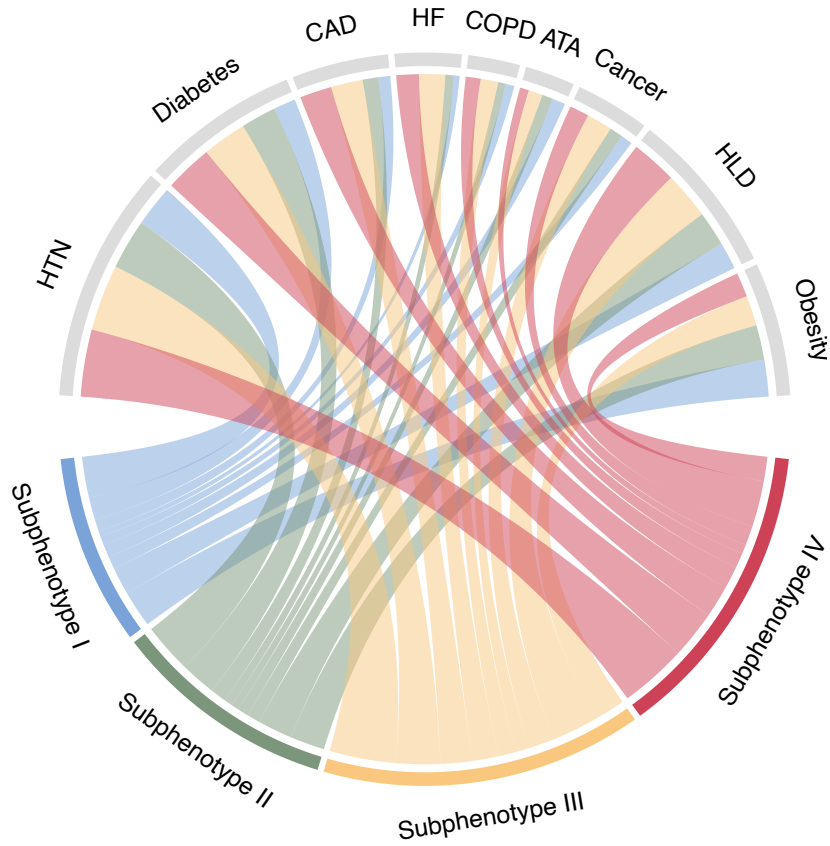
Abnormal biomarkers vs. Subphenotype III



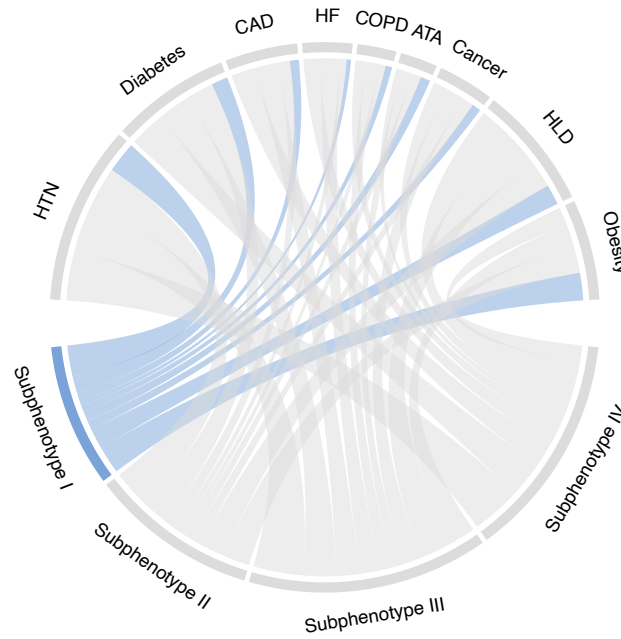
Abnormal biomarkers vs. Subphenotype IV



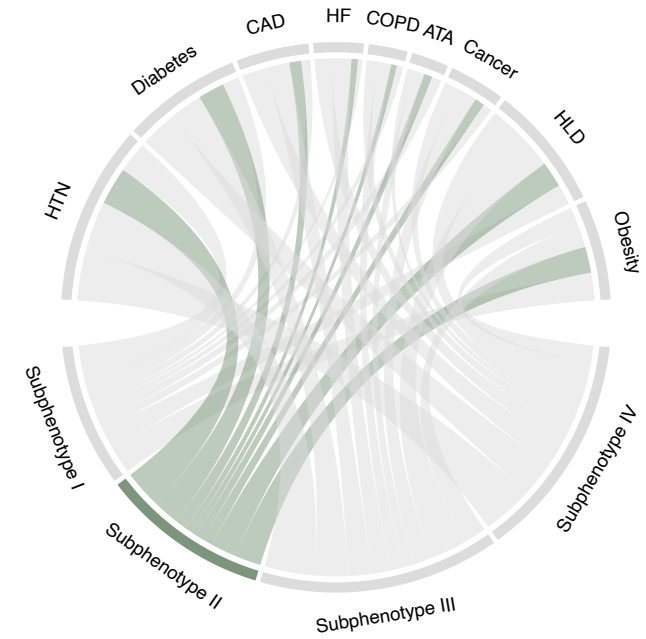
Comorbidity vs. subphenotypes



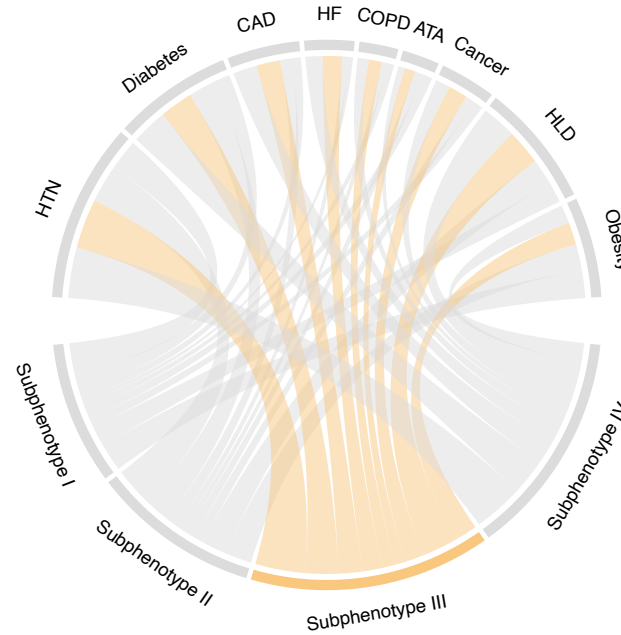
Comorbidity vs. Subphenotype I



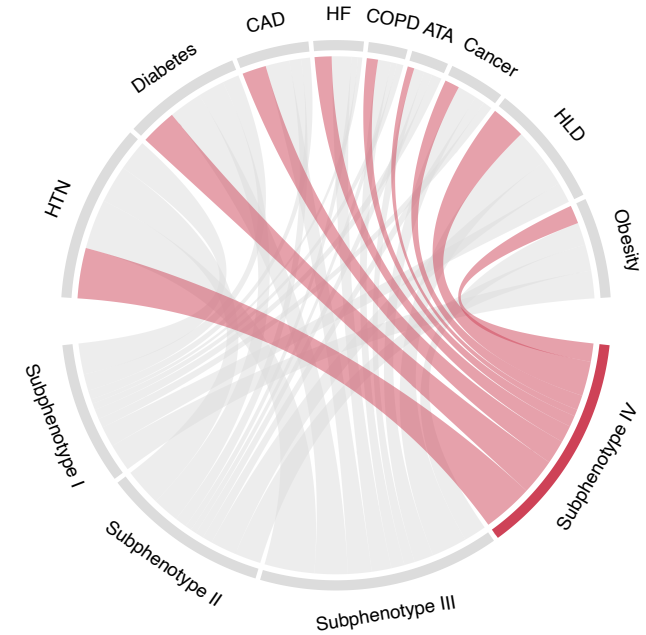
Comorbidity vs. Subphenotype II

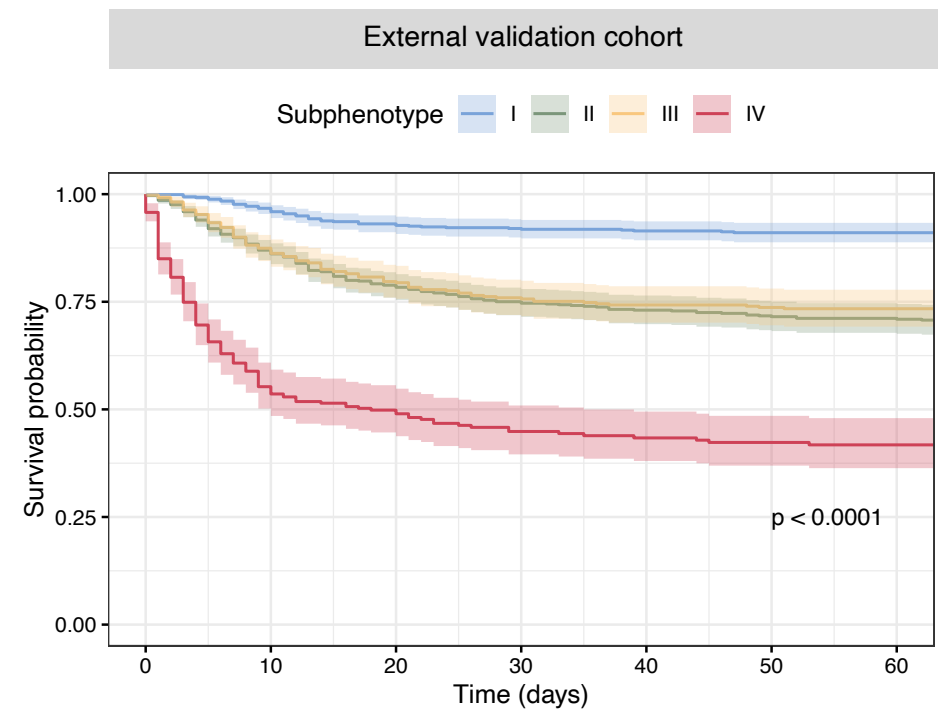
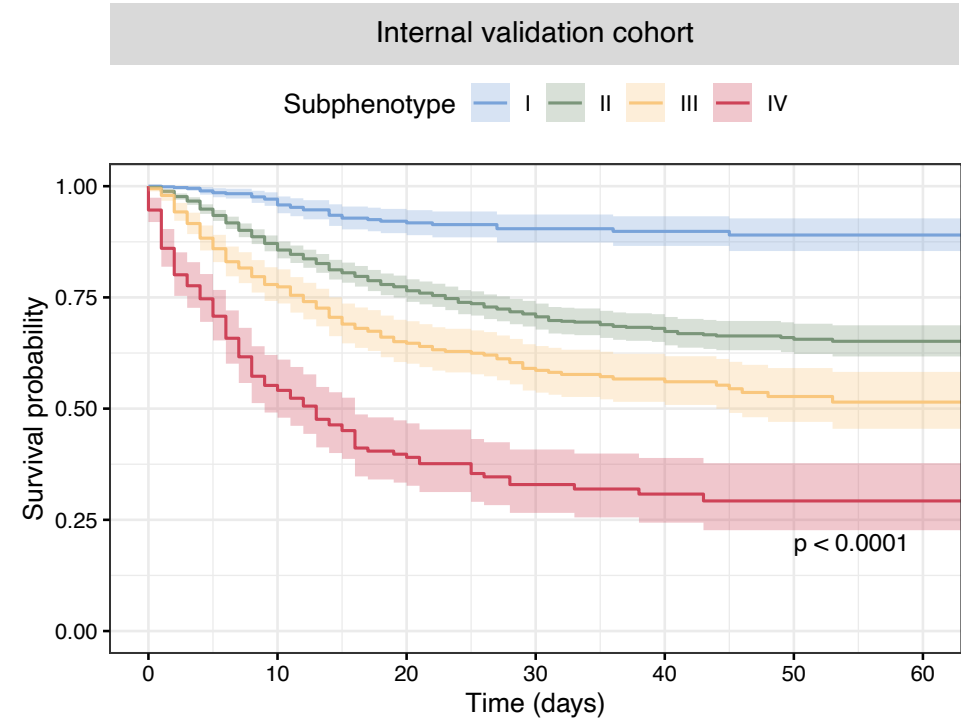
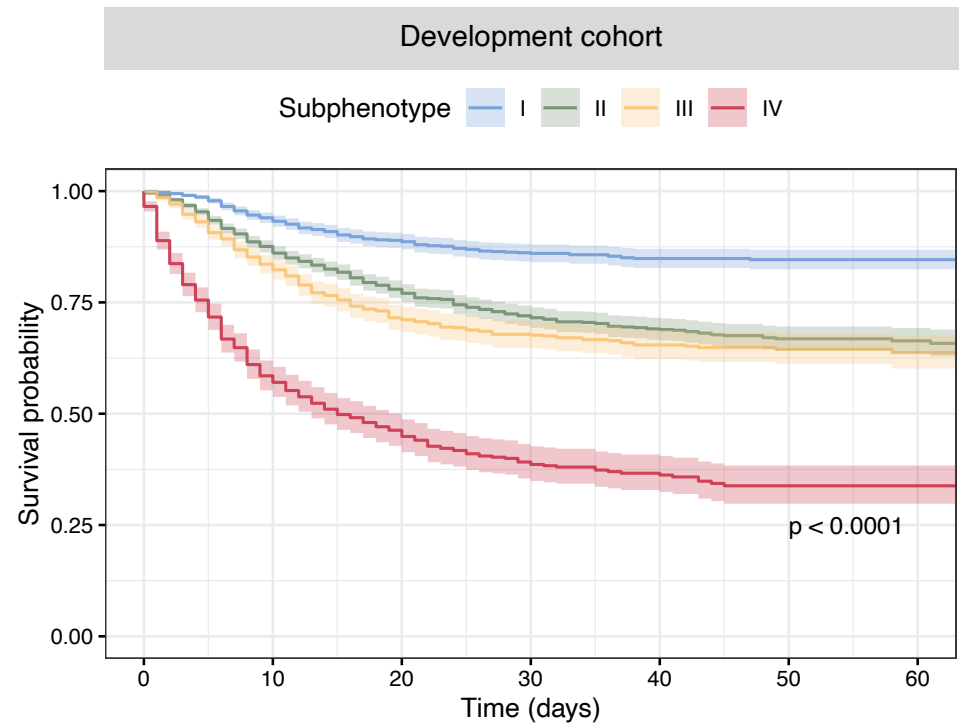


Comorbidity vs. Subphenotype III



Comorbidity vs. Subphenotype IV

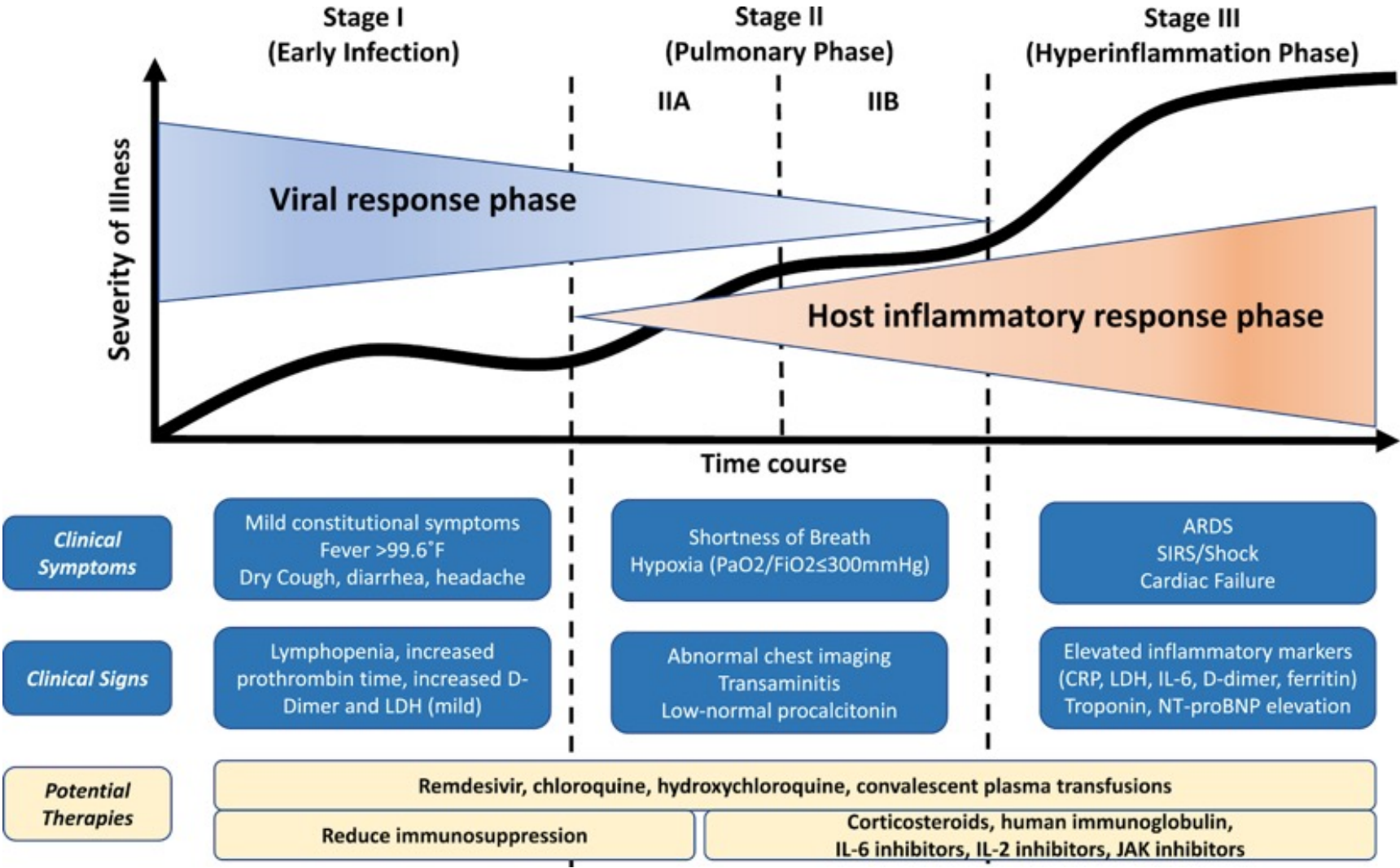




Outline

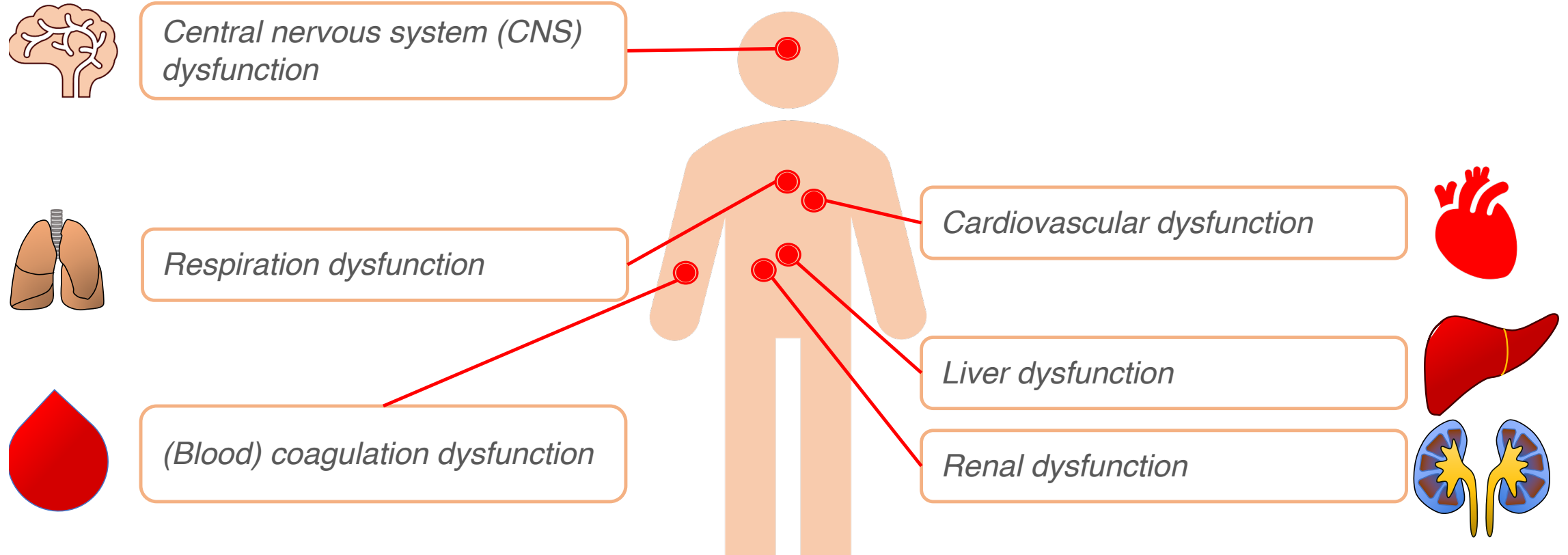
- Introduction
- Subphenotyping of COVID-19 at infection confirmation
- **Subphenotyping of Severe COVID-19 after Mechanical Ventilation**
- Subphenotyping of Long COVID
- Subphenotyping of PD
- Discussions

An Early Conceptual Model for the Progression of COVID-19 in Acute Phase



Siddiqi, Hasan K., and Mandeep R. Mehra. "COVID-19 illness in native and immunosuppressed states: A clinical–therapeutic staging proposal." *The Journal of Heart and Lung Transplantation* 39, no. 5 (2020): 405.

Sequential Organ Failure Assessment



Cohort Setup

Mar 1st to May 12th

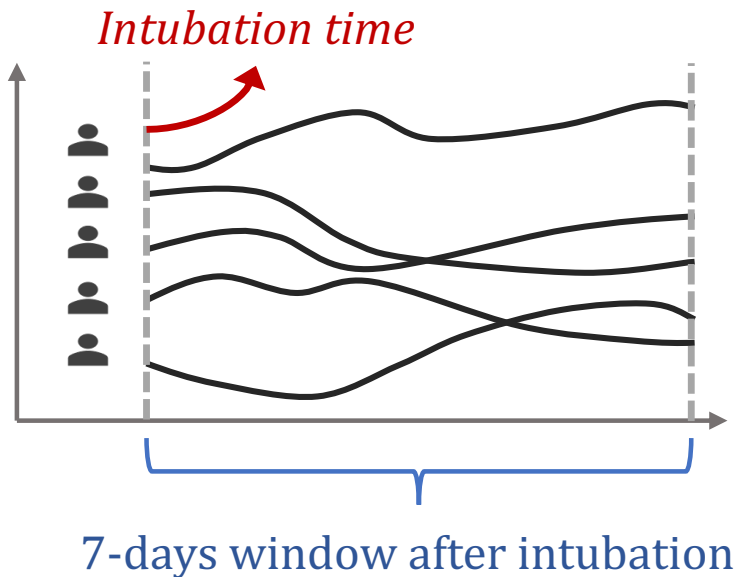
*New York Presbyterian
Weill Cornell Medical Center
(NYP-WCMC) – **348 patients***

Development

*New York Presbyterian
Lower Manhattan Hospital
(NYP-LMH) – **100 patients***

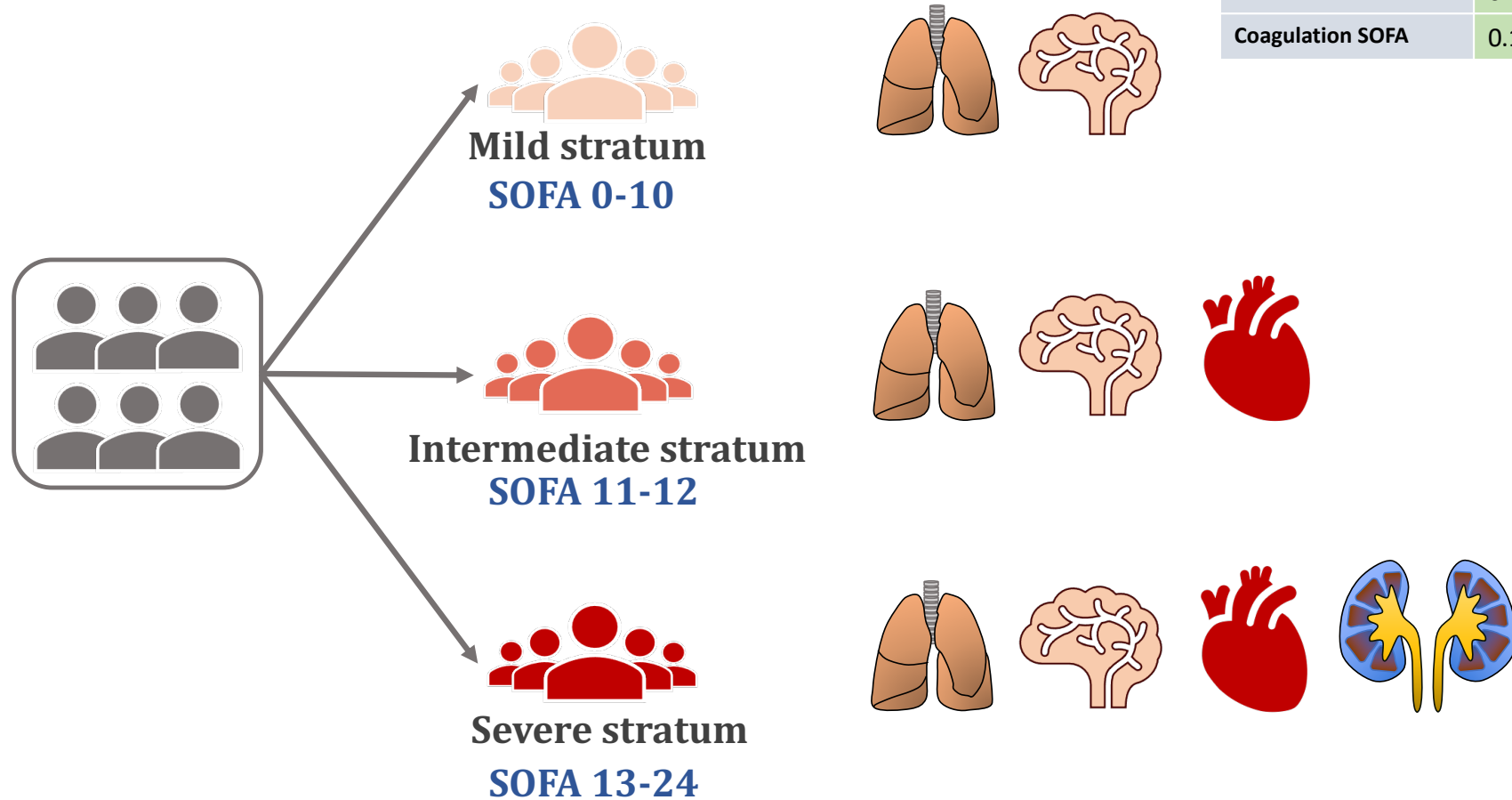
Validation

Post-intubation SOFA trajectory



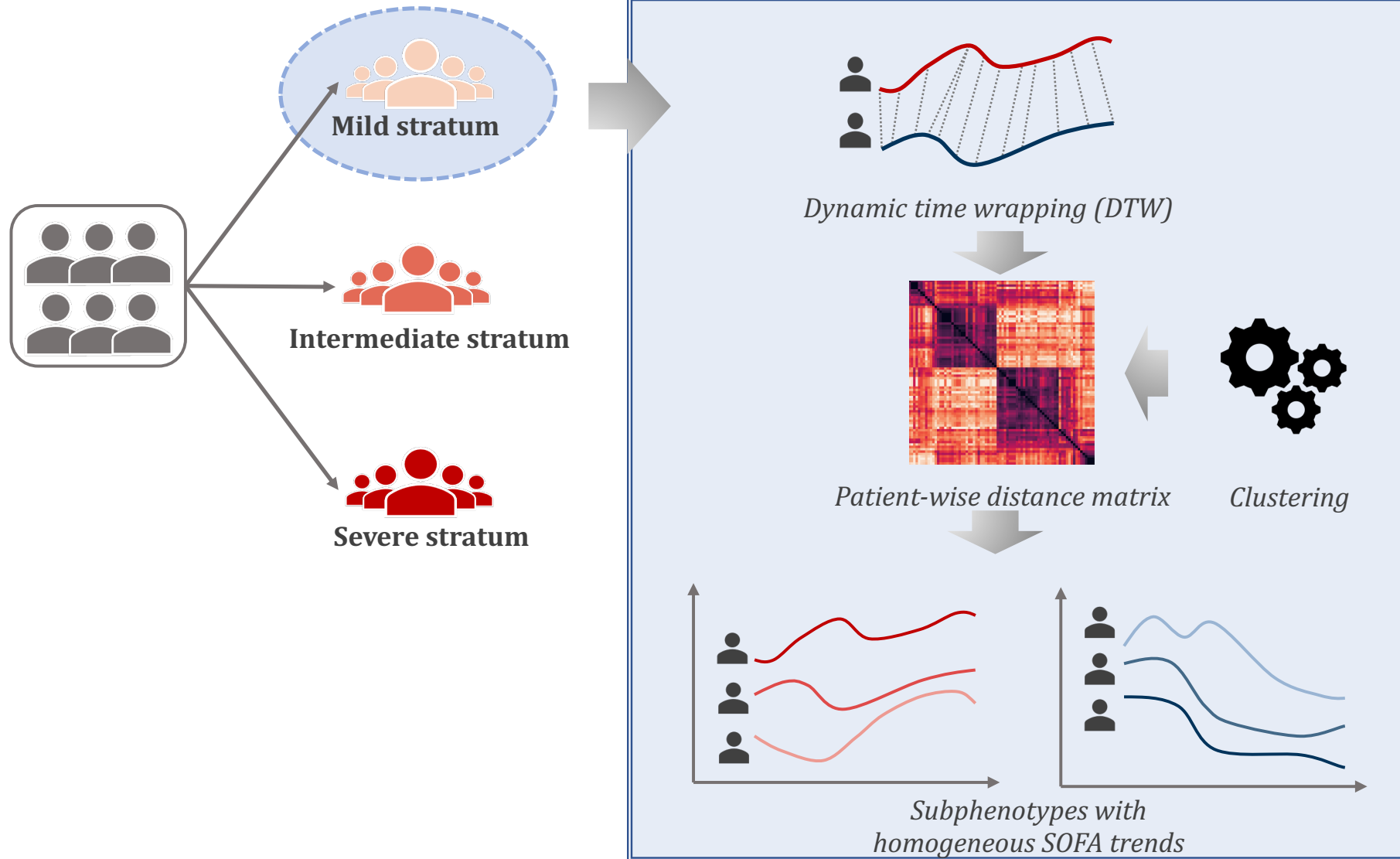
- SOFA scores were calculated every 24 hours
- Patients missing more than 3-days SOFA data were excluded

Baseline Stratification

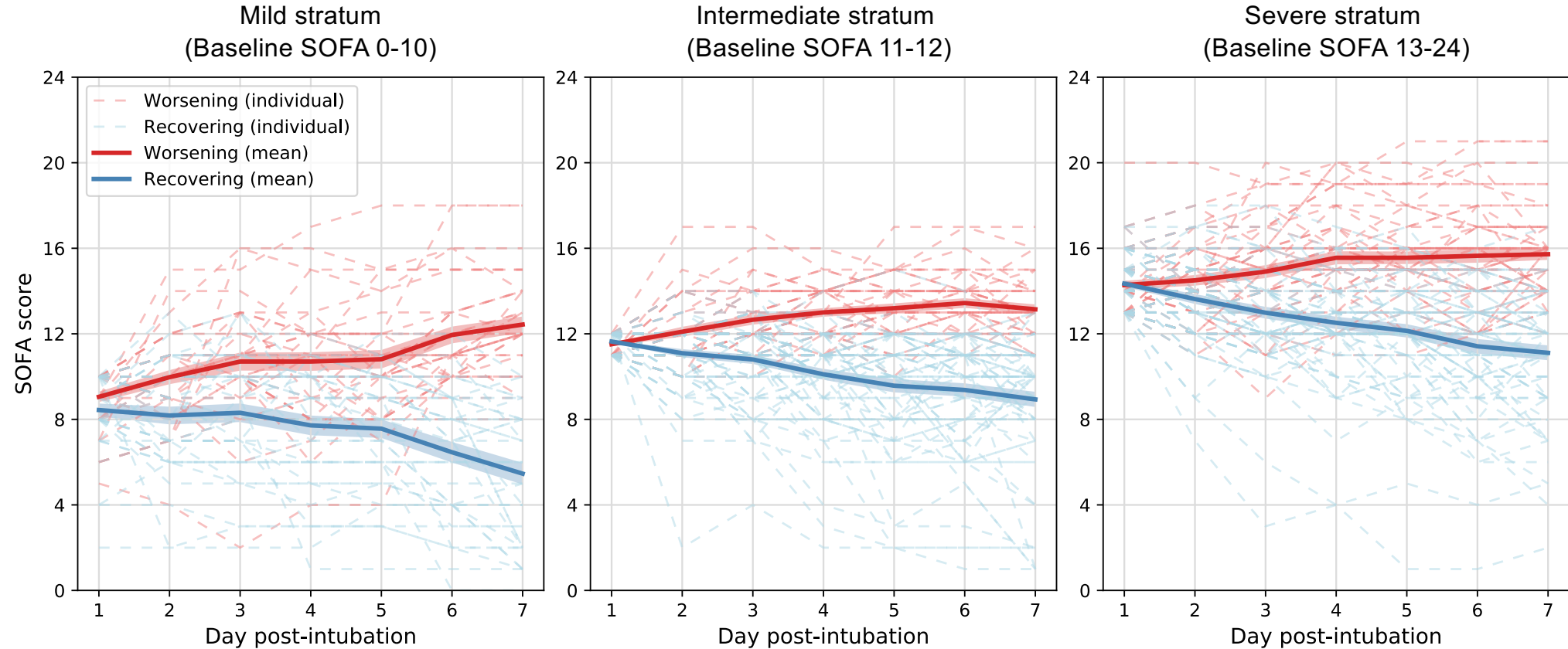


	Mild	Intermediate	Severe
Respiration SOFA	3.45 (0.89)	3.89 (0.45)	3.97 (0.25)
CNS SOFA	3.34 (1.13)	3.72 (0.47)	3.94 (0.24)
Cardiovascular SOFA	1.32 (1.34)	3.41 (0.88)	3.69 (0.70)
Renal SOFA	0.16 (0.54)	0.35 (0.67)	1.96 (1.44)
Liver SOFA	0.20 (0.46)	0.14 (0.43)	0.37 (0.67)
Coagulation SOFA	0.12 (0.40)	0.04 (0.20)	0.28 (0.64)

Trajectory Grouping



Identified Subphenotypes

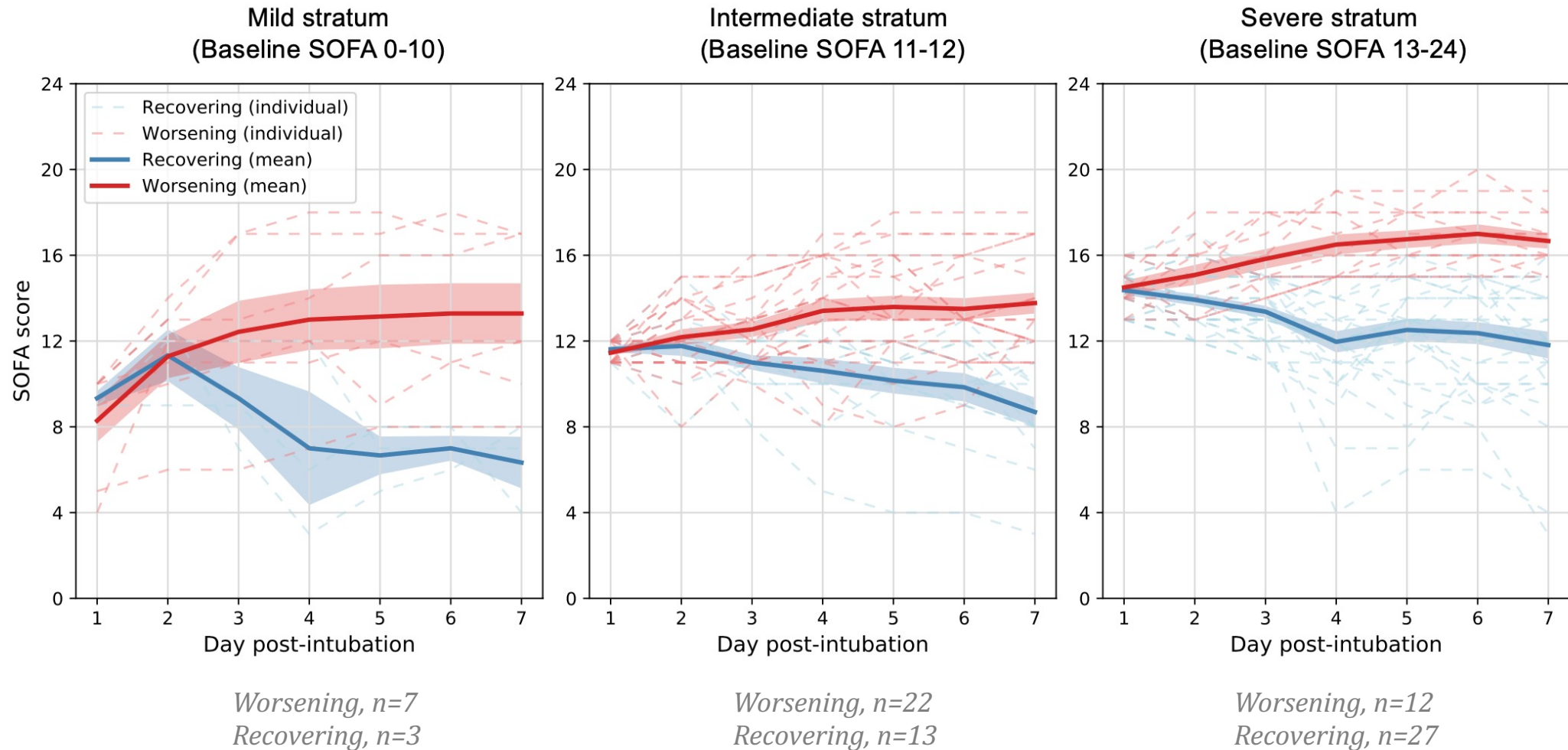


Worsening, n=37
Recovering, n=39

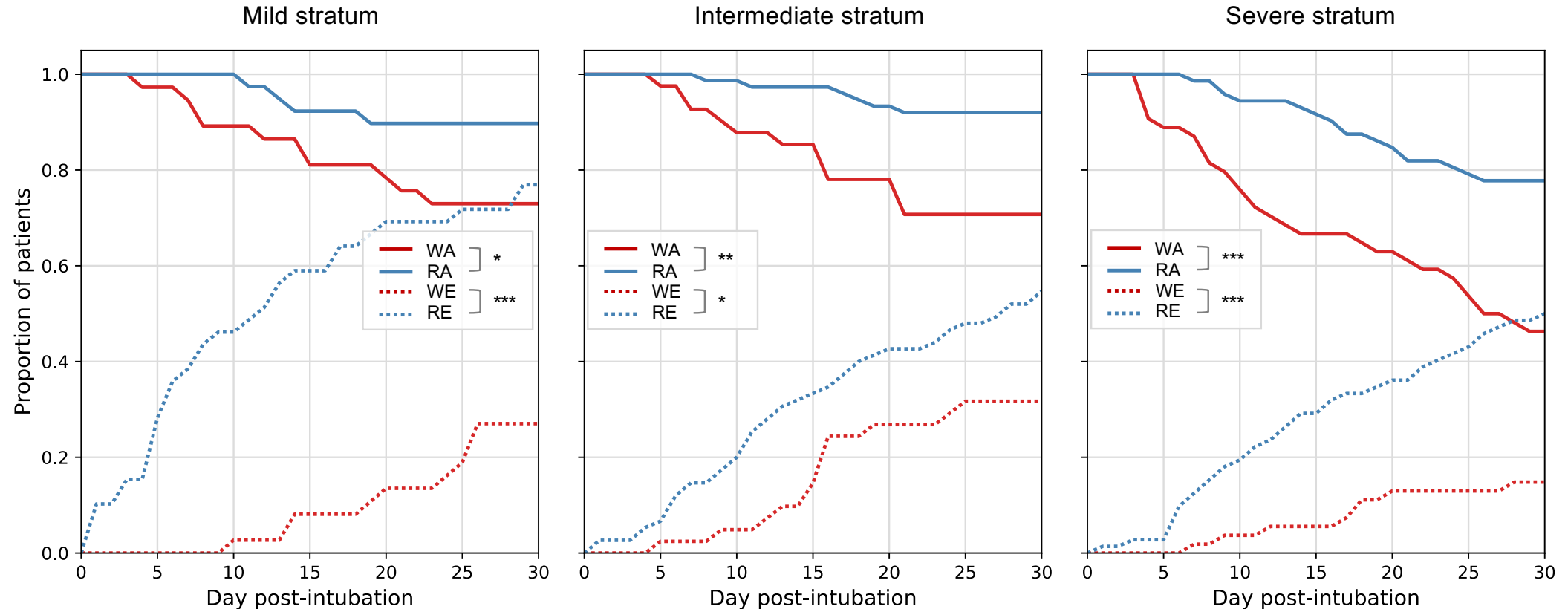
Worsening, n=41
Recovering, n=75

Worsening, n=54
Recovering, n=72

Identified Subphenotypes



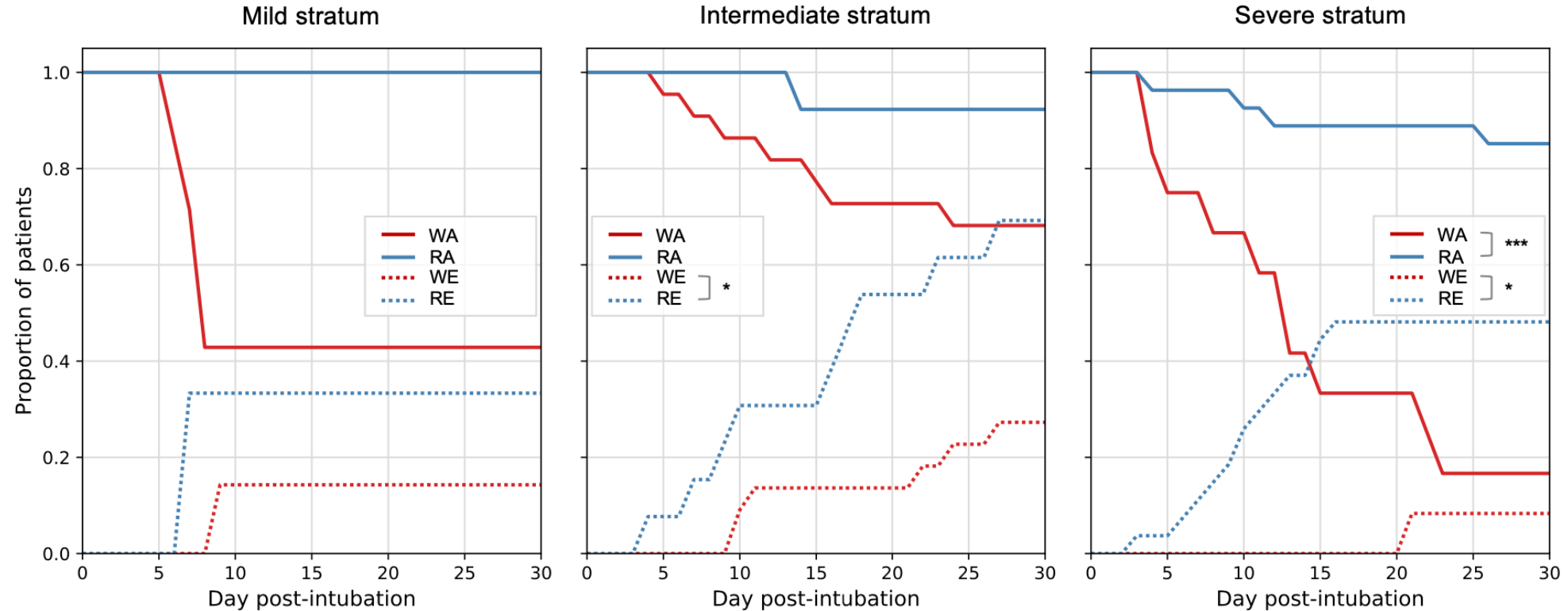
Association with Outcomes



* $P < 0.05$
 ** $p < 0.01$
 *** $p < 0.001$

Abbreviations: WA=worsening subphenotype alive; RA=recovering subphenotype alive;
 WE=worsening subphenotype extubated; RE=recovering subphenotype extubated.

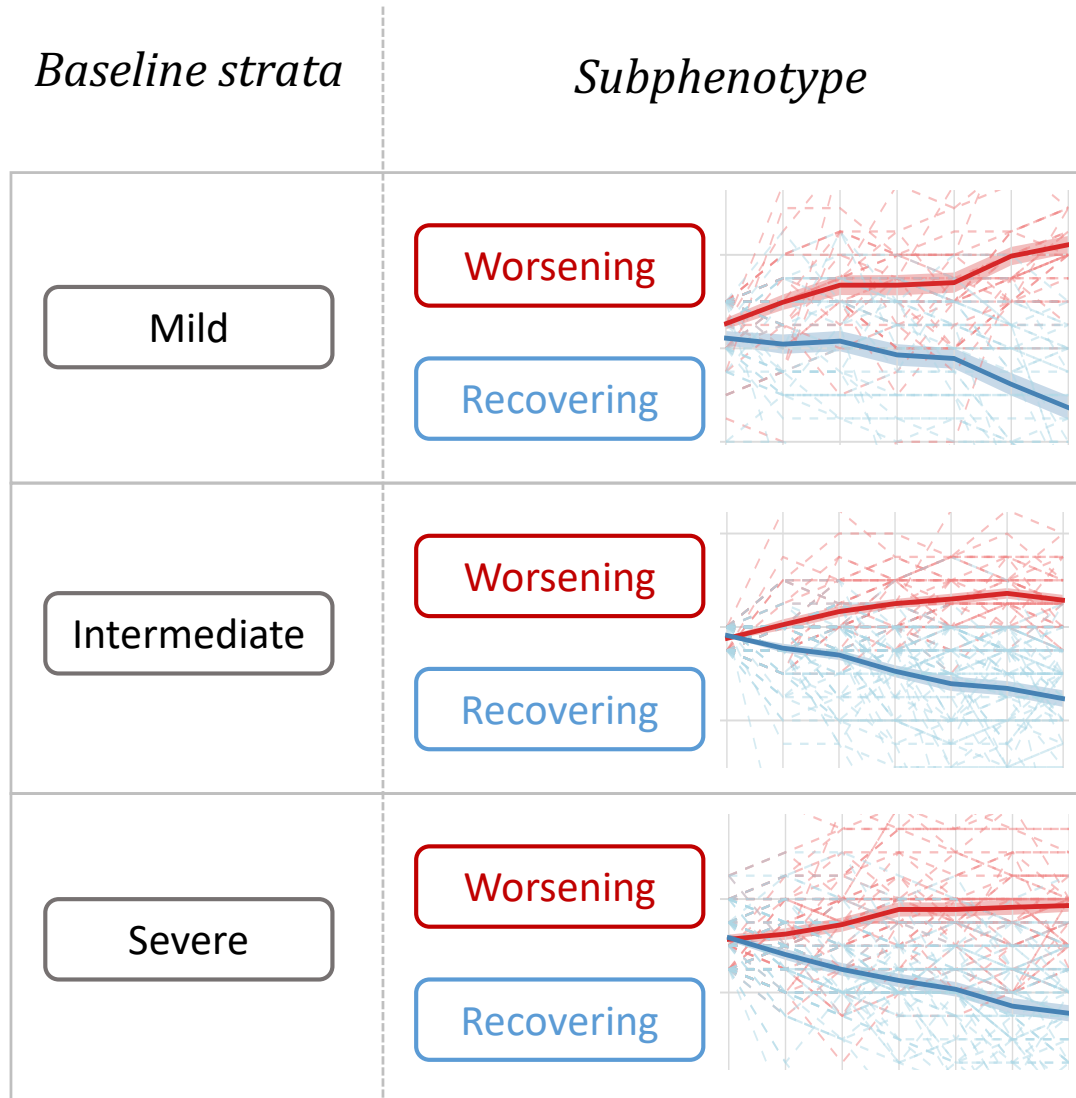
Association with Outcomes



* $P < 0.05$
 ** $p < 0.01$
 *** $p < 0.001$

Abbreviations: WA=worsening subphenotype alive; RA=recovering subphenotype alive;
 WE=worsening subphenotype extubated; RE=recovering subphenotype extubated.

Discriminative Biomarker Identification



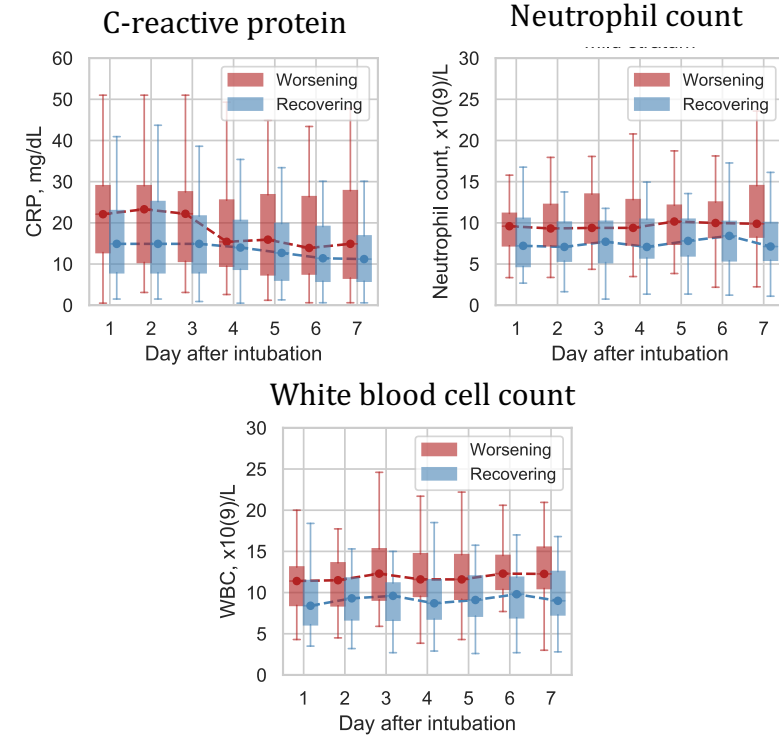
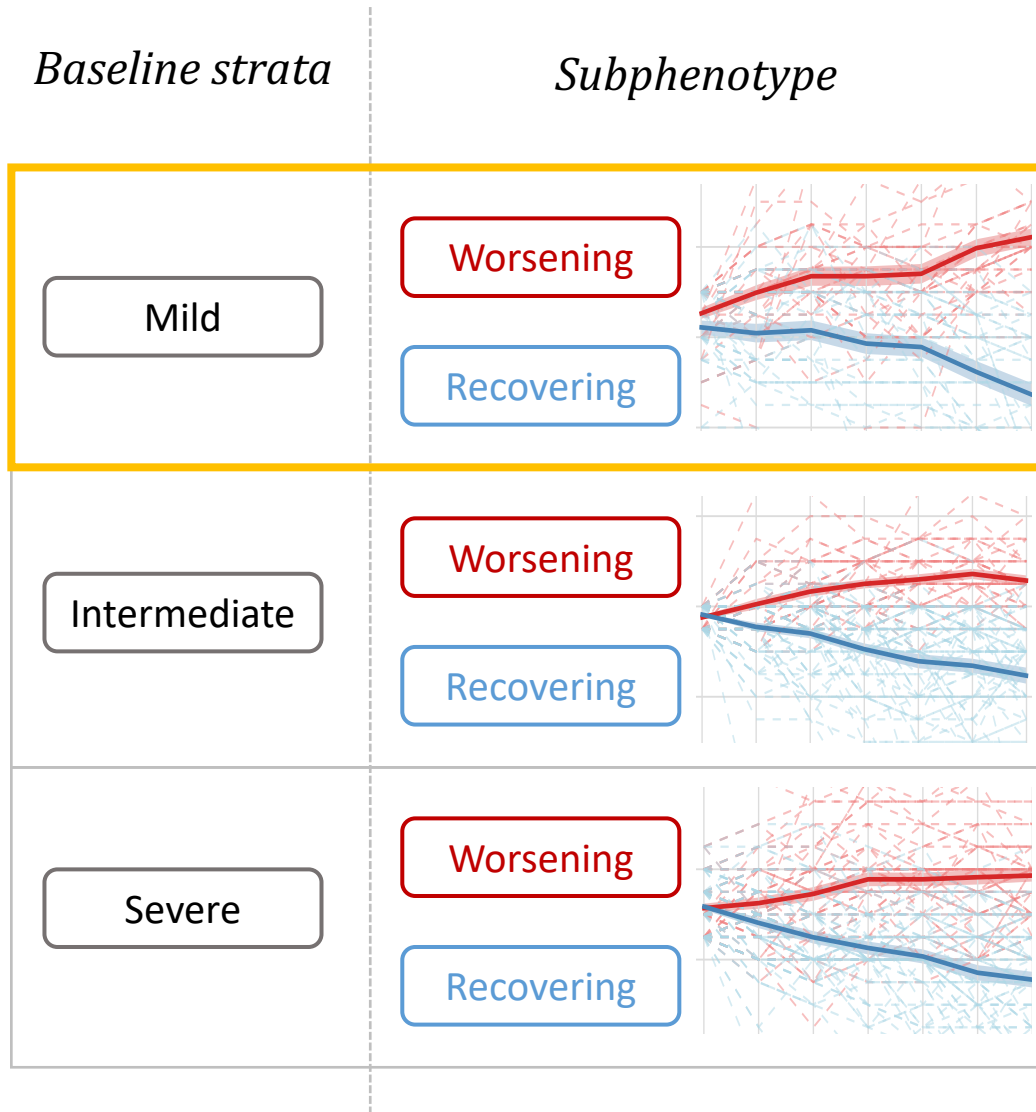
- 27 laboratory test values

E.g., albumin, bilirubin, and white blood cell count

- 4 vital signs

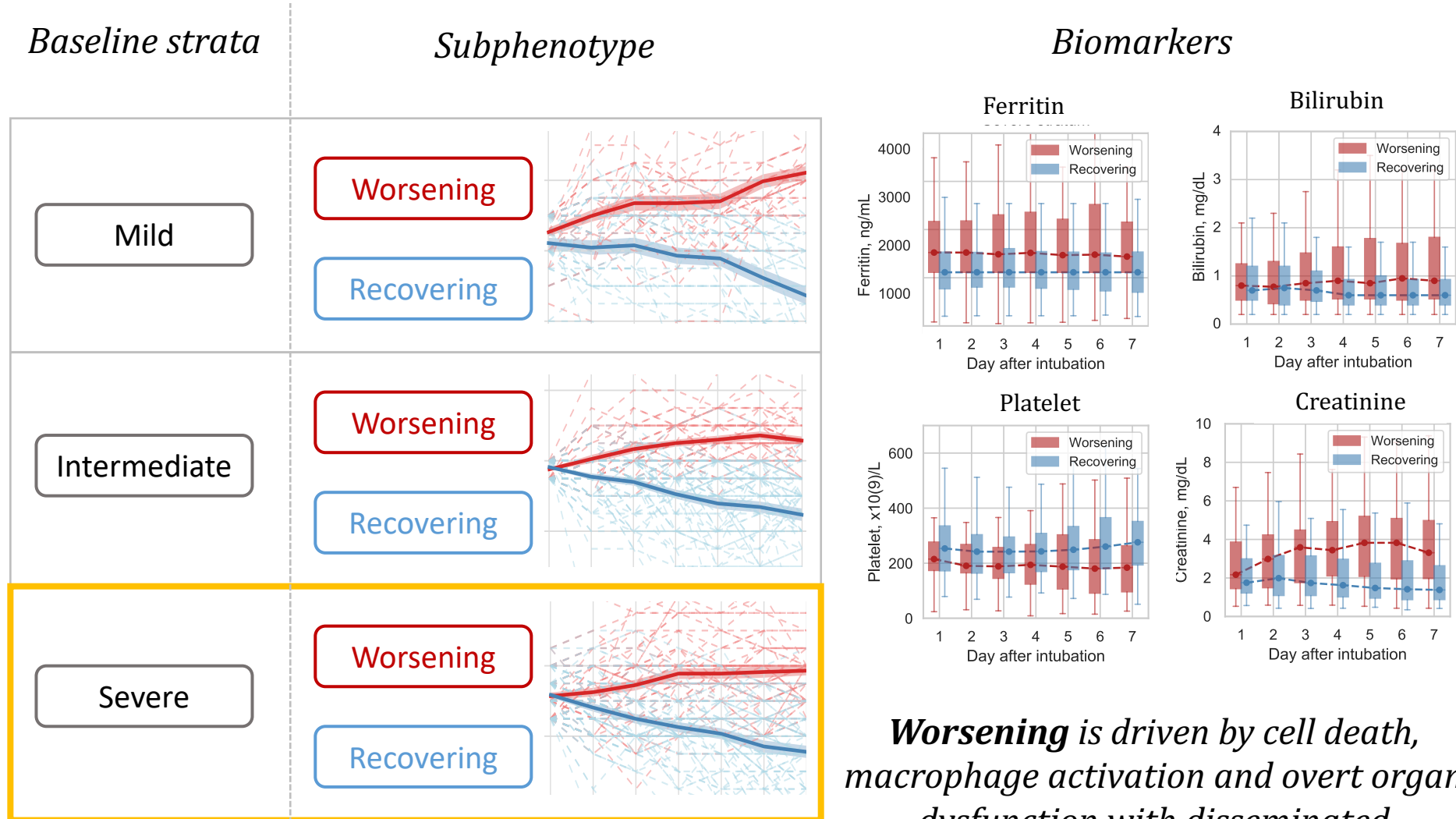
Body temperature, urine output, mean arterial pressure, and Glasgow Coma Scale (GCS, consciousness)

Discriminative Biomarker Identification



Inflammatory hinders the worsening subphenotype from recovery.

Discriminative Biomarker Identification



***Worsening** is driven by cell death, macrophage activation and overt organ dysfunction with disseminated intravascular coagulation*

Outline

- Introduction
- Subphenotyping of COVID-19 at infection confirmation
- Subphenotyping of Severe COVID-19 after Mechanical Ventilation
- **Subphenotyping of Long COVID**
- Subphenotyping of PD
- Discussions

Pipeline

DATABASE



Electronic Health Records (EHR) data for patients with lab-confirmed SARS-CoV-2 Infection from two clinical research networks (CRN)

- INSIGHT: New York City area
- OneFlorida+: Florida, Georgia and Alabama

METHOD

Cohort: SARS-CoV-2 infected patients with newly incident conditions within 30-180 days after infection

Variables: 137 investigative conditions

Anemia	1	0	0	0
Circulatory problem	1	0	0	0
Disorders of stomach	0	1	0	0
Malaise and fatigue	0	1	1	1
Nausea and Vomiting	1	1	0	...
Headache	0	0	1	0
Respirator problem	0	0	1	0
...
Musculoskeletal pain	0	0	0	1
Osteoarthritis	0	0	0	1

Step 1. Binary vector representations of patients with incident PASC diagnosis



Topic 1

PASC	Weight
Disorders of stomach	0.32
Nausea and Vomiting	0.23
Esophageal disorders	0.20
...	...

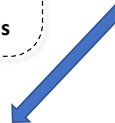
Topic 2

PASC	Weight
Anemia	0.36
Heart failure	0.25
Cardiac dysrhythmias	0.19
...	...

Topic K

PASC	Weight
Musculoskeletal pain	0.34
Osteoarthritis	0.28
Spondylopathies	0.14
...	...

Step 2. Inference of PASC topics. Each PASC topic encodes a particular set of frequently co-occurred incident PASC diagnosis



Topic 1	0.05	0.74	0.08	0.02
Topic 2	0.65	0.01	0.13	0.06
...
Topic K-1	0.12	0.03	0.53	0.03
Topic K	0.07	0.10	0.02	0.69

Step 3. Derivation of the patient representation in the PASC topic space



Subphenotype 1

Subphenotype 2

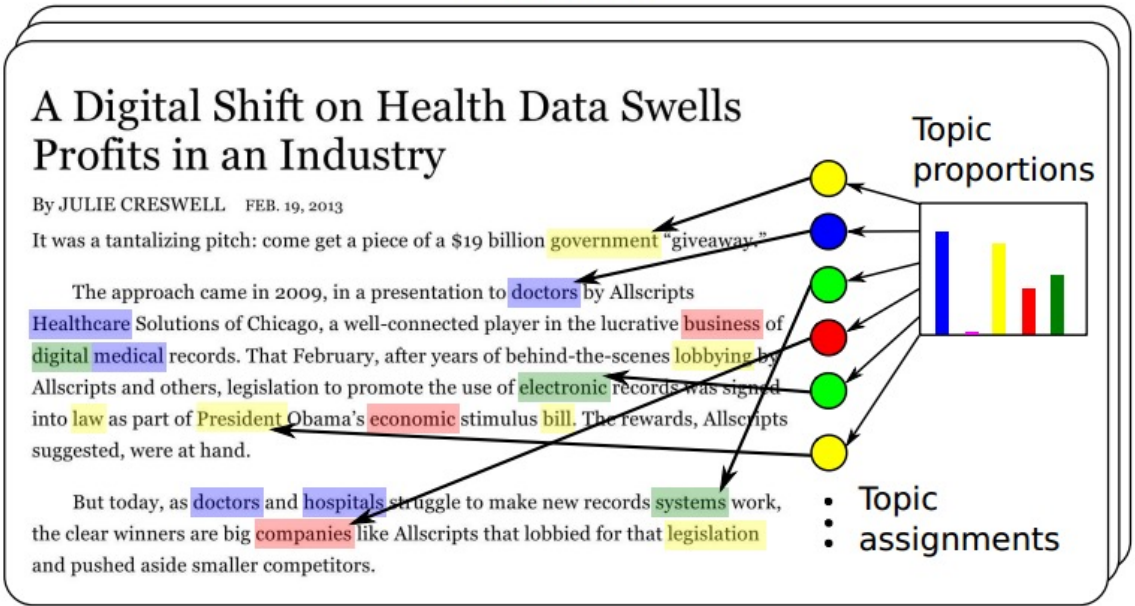
Subphenotype 3

Subphenotype 4

Step 4. Derivation of the PASC subphenotypes as patient groups with the PASC topic-based representation through cluster analysis.

Background – probabilistic topic modeling

Documents



- **Topics:** a distribution over all words. Global parameters shared by all documents.
- **Topic proportion:** represent the importance of each topic in representing this document. Local parameters denotes the new feature of each document in topic space.

Topics

health 0.03	team 0.03	government 0.04	business 0.04	computer 0.03
medical 0.03	basketball 0.02	law 0.02	money 0.02	system 0.02
disease 0.02	points 0.01	politics 0.01	economic 0.02	software 0.02
hospital 0.01	score 0.01	legislation 0.01	company 0.01	program 0.01
...

Topic Modeling

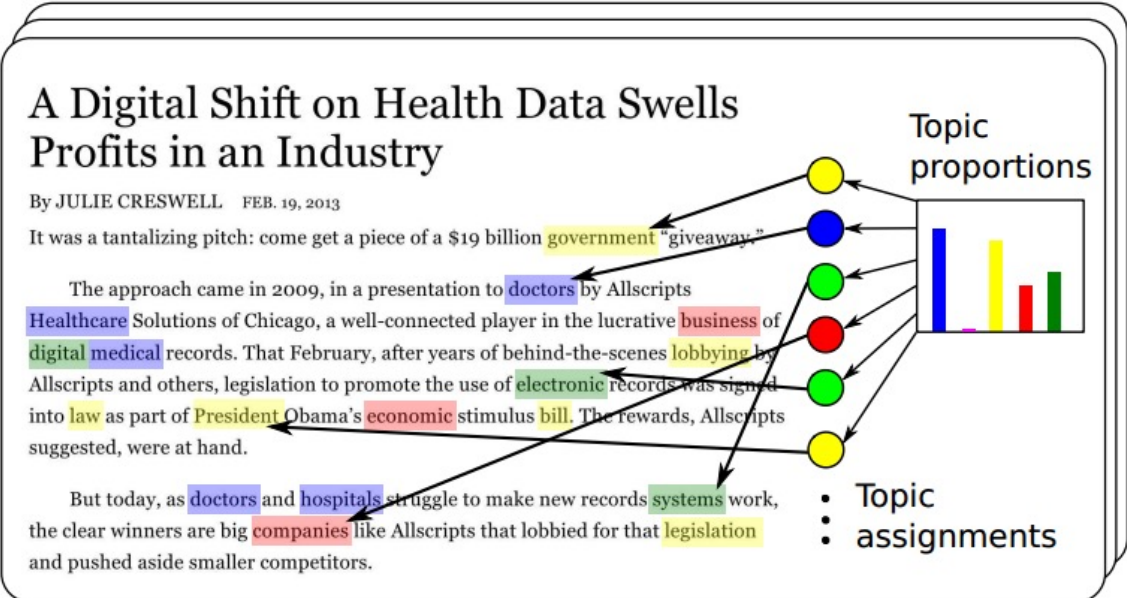
Transform the document represented by words to a new feature represented by topics.

The basic idea of topic modeling: exploring different topics (group of some “similar” words) for documents

Similar (word cooccurrence): some words should often appear simultaneously in one document

Background – probabilistic topic modeling

Documents



Topics

health 0.03	team 0.03	government 0.04	business 0.04	computer 0.03
medical 0.03	basketball 0.02	law 0.02	money 0.02	system 0.02
disease 0.02	points 0.01	politics 0.01	economic 0.02	software 0.02
hospital 0.01	score 0.01	legislation 0.01	company 0.01	program 0.01
...

Topic Modeling

Transform the document represented by words to a new feature represented by topics.

Vocabulary	1)	2)
industry	3	0
NBA	0	6
economy	2	1
economic	3	2
game	0	5
medical	2	0
business	3	4
...



VS

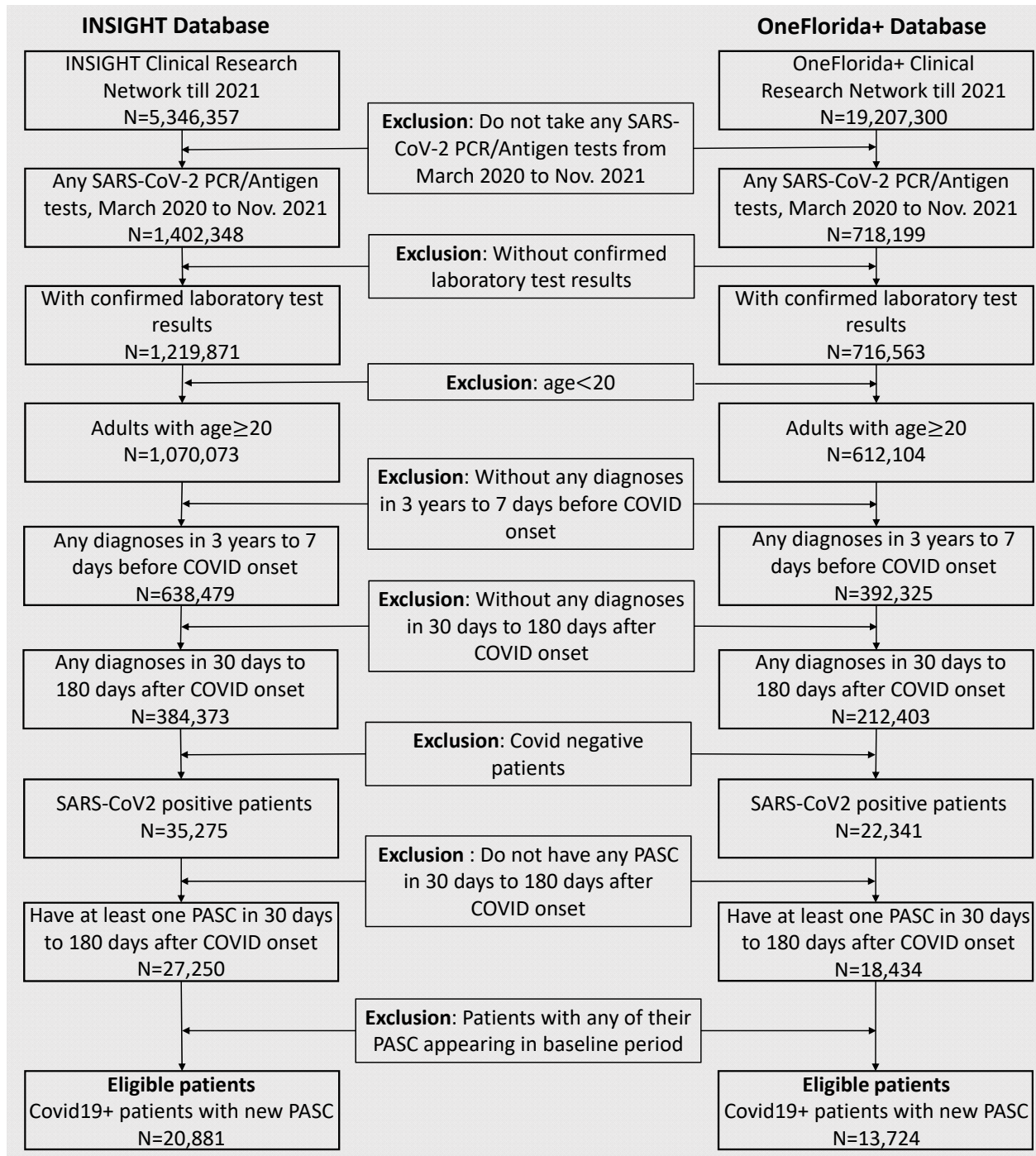
Vocabulary	1)	2)
Topic1	0.4	0.05
Topic2	0.06	0.3
Topic3	0.3	0.04
Topic4	0.05	0.3
Topic5	0.19	0.31

... This vector should be high-dimensional, sparse, and discrete.

This vector should be low-dimensional, dense, and continuous

Topic modeling	PCA
Topics	Factors
Topic proportion	Factor loading

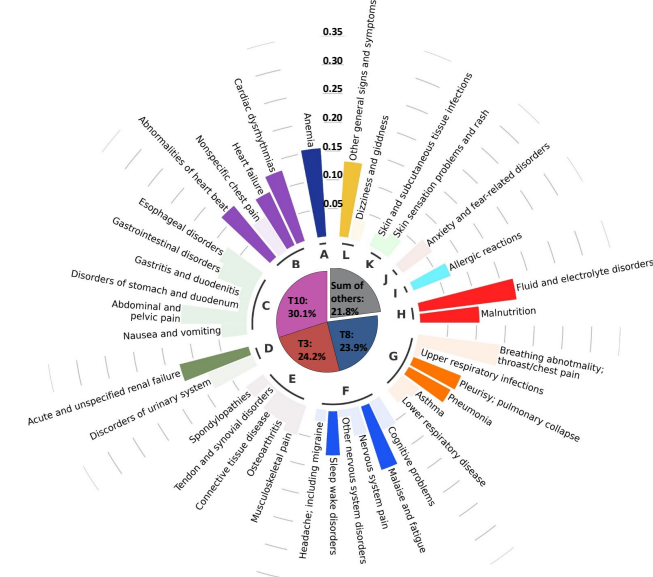
Inclusion-Exclusion cascade



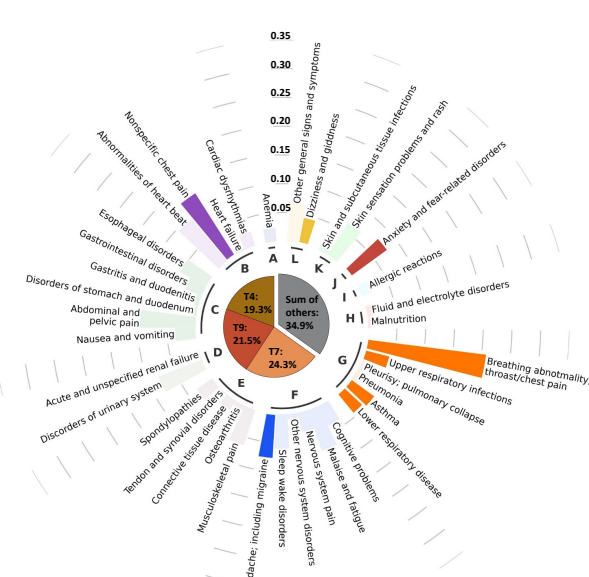
INSIGHT

CCSR domain	PASC	PASC Topic																			
		T1	T2	T3	T4	T5	T6	T7	T8	T9	T10										
Diseases of the Blood	Anemia																				
Diseases of the Circulatory System	Cardiac dysrhythmias																				
	Heart failure																				
	Nonspecific chest pain																				
	Abnormalities of heart beat																				
Diseases of the Digestive System	Esophageal disorders																				
	Gastrointestinal disorders																				
	Gastritis and duodenitis																				
	Disorders of stomach and duodenum																				
	Abdominal and pelvic pain																				
	Nausea and vomiting																				
Diseases of the Genitourinary System	Acute and unspecified renal failure																				
Diseases of the Musculoskeletal System and Connective Tissue	Spondylopathies																				
	Tendon and synovial disorders																				
	Connective tissue disease																				
	Osteoarthritis																				
Diseases of the Nervous System	Musculoskeletal pain																				
	Headache																				
	Sleep wake disorders																				
	Other nervous system disorders																				
	Nervous system pain																				
	Malaise and fatigue																				
Diseases of the Respiratory System	Cognitive problems																				
	Lower respiratory disease																				
Endocrine, Nutritional and Metabolic Diseases	Asthma																				
	Breathing abnormality and throat/chest pain																				
Fluid and electrolyte disorders																					
Anxiety and fear-related disorders																					
Other general signs and symptoms																					
Diseases of the Skin and Subcutaneous Tissue	Skin sensation problems and rash																				
	Skin and subcutaneous tissue infections																				
Other PASC (106 items)																					

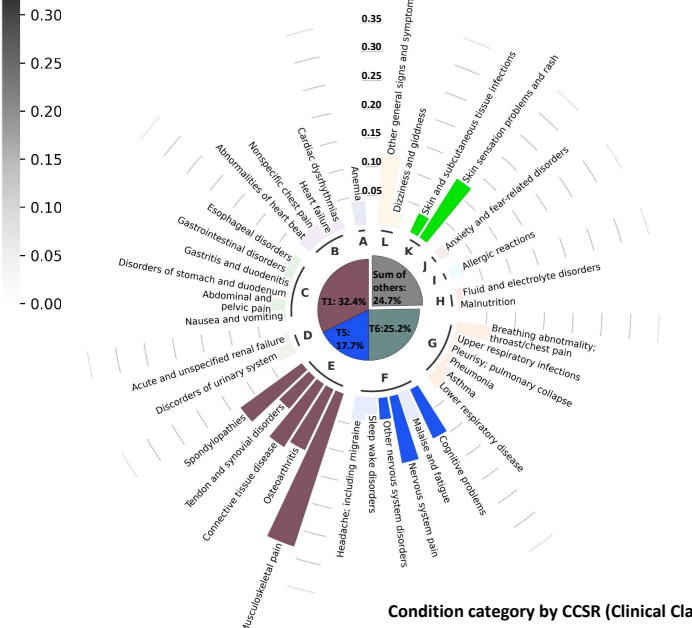
Subphenotype 1 (Cardiac and Renal)



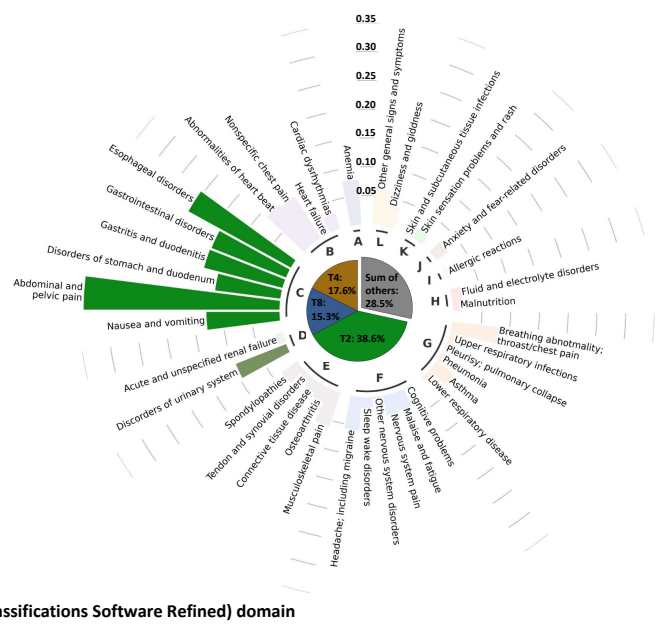
Subphenotype 2 (Respiratory, Sleep and Anxiety)



Subphenotype 3 (Musculoskeletal and Nervous)



Subphenotype 4 (Digestive and Respiratory)

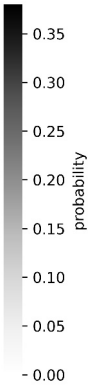


Condition category by CCSR (Clinical Classifications Software Refined) domain

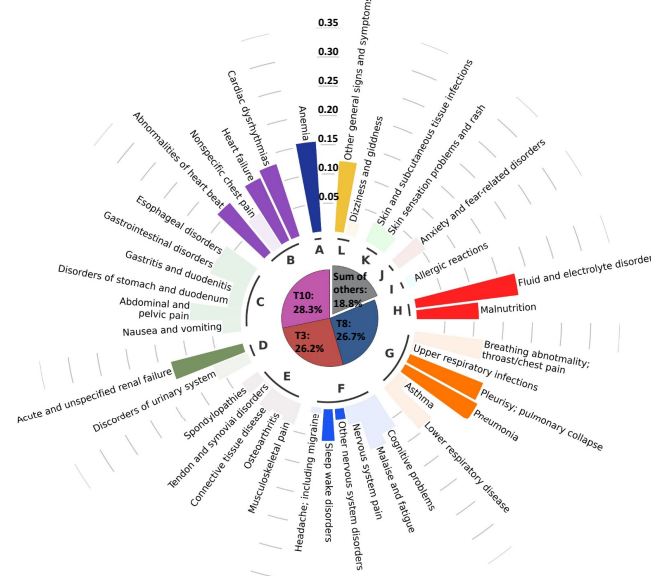
- A: Diseases of the blood and blood-forming organs
- B: Diseases of the circulatory system
- C: Diseases of the digestive system
- D: Diseases of the genitourinary system
- E: Diseases of the musculoskeletal system
- F: Diseases of the nervous system
- G: Diseases of the respiratory system
- H: Endocrine, nutritional and metabolic diseases
- I: Injury and poisoning
- J: Mental and behavioral disorders
- K: Diseases of the skin and subcutaneous tissue
- L: Others

OneFlorida+

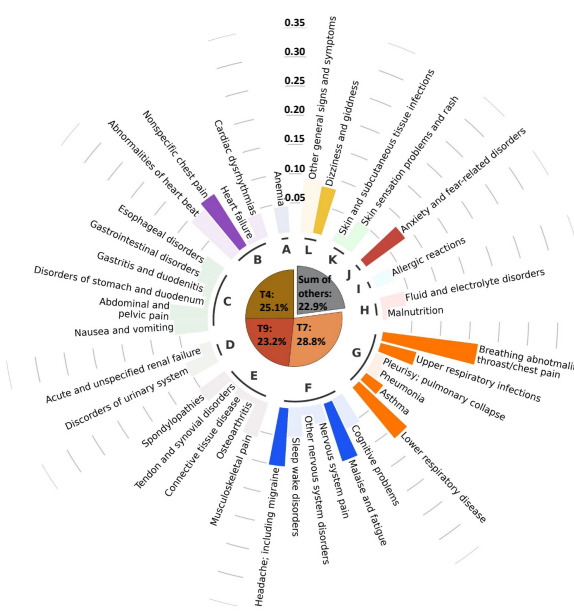
CCSR domain	PASC	PASC Topic									
		T1	T2	T3	T4	T5	T6	T7	T8	T9	T10
Diseases of the Blood	Anemia										
Diseases of the Circulatory System	Cardiac dysrhythmias										
	Heart failure										
	Nonspecific chest pain										
	Abnormalities of heart beat										
Diseases of the Digestive System	Esophageal disorders										
	Gastrointestinal disorders										
	Gastritis and duodenitis										
	Disorder of stomach and duodenum										
	Abdominal and pelvic pain										
	Nausea and vomiting										
Diseases of the Genitourinary System	Acute and unspecified renal failure										
Diseases of the Musculoskeletal System and Connective Tissue	Spondylopathies										
	Tendon and synovial disorders										
	Connective tissue disease										
	Osteoarthritis										
Diseases of the Nervous System	Musculoskeletal pain										
	Headache										
	Sleep wake disorders										
	Other nervous system disorders										
	Nervous system pain										
	Malaise and fatigue										
Diseases of the Respiratory System	Cognitive problems										
	Lower respiratory disease										
	Asthma										
	Pneumonia										
	Pleurisy; Pulmonary collapse										
Upper respiratory disease											
Breathing abnormality and throat/chest pain											
Endocrine, Nutritional and Metabolic Diseases	Fluid and electrolyte disorders										
Neurological Disorders	Anxiety and fear-related disorders										
Other Symptoms and Signs	Other general signs and symptoms										
Diseases of the Skin and Subcutaneous Tissue	Skin sensation problems and rash										
	Skin and subcutaneous tissue infections										
	Other PASC (103 items)										



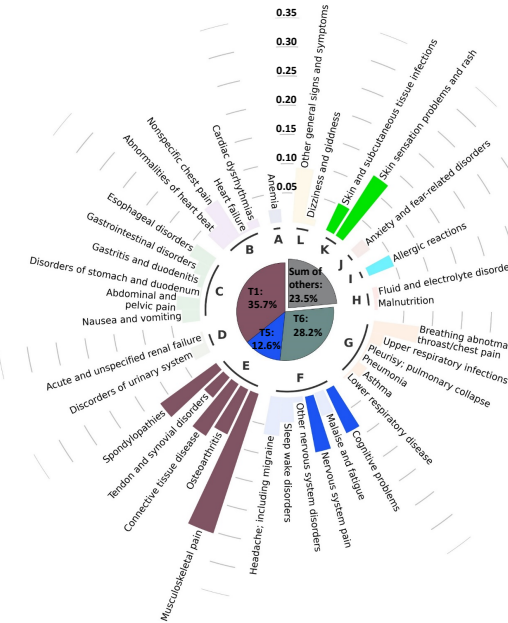
Subphenotype 1 (Cardiac and Renal)



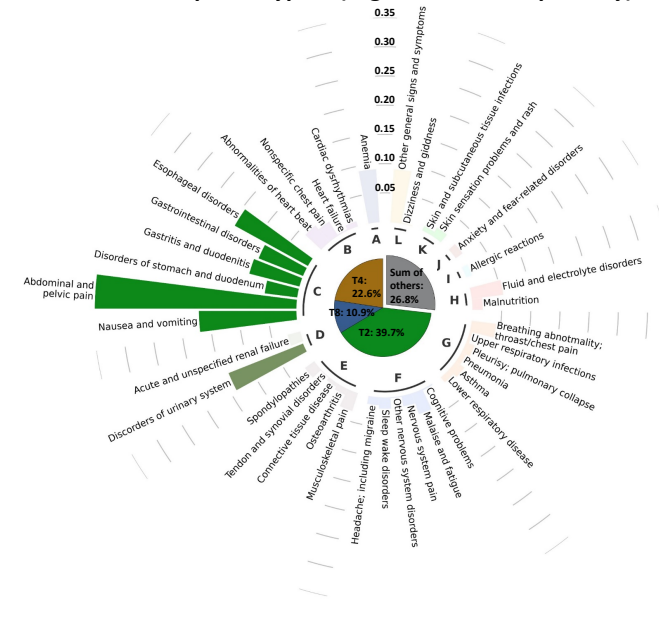
Subphenotype 2 (Respiratory, Sleep and Anxiety)



Subphenotype 3 (Musculoskeletal and Nervous)



Subphenotype 4 (Digestive and Respiratory)

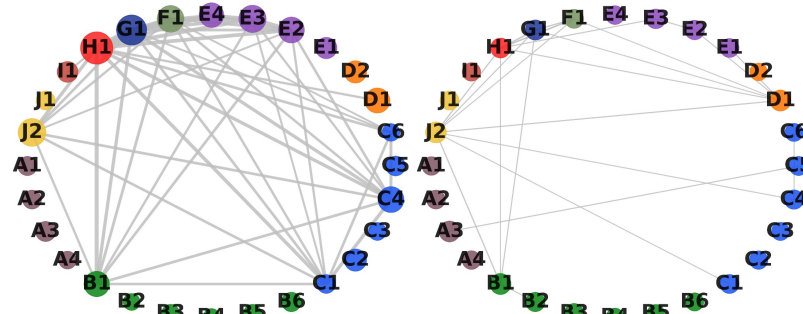


PASC category by CCSR domain

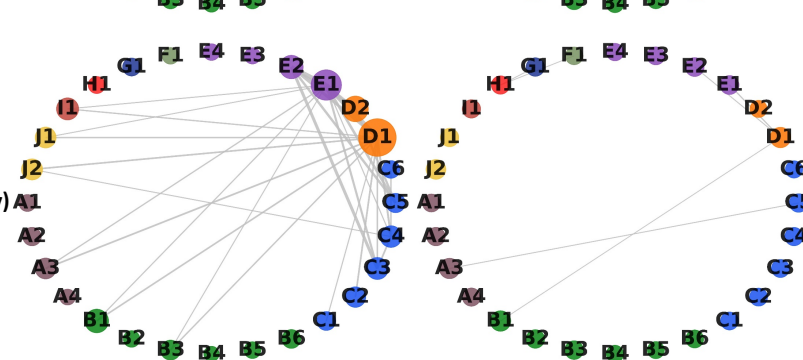
- A: Diseases of the blood and blood-forming organs
- B: Diseases of the circulatory system
- C: Diseases of the digestive system
- D: Diseases of the genitourinary system
- E: Diseases of the musculoskeletal system
- F: Diseases of the nervous system
- G: Diseases of the respiratory system
- H: Endocrine, nutritional and metabolic diseases
- I: Injury and poisoning
- J: Diseases of the nervous system
- K: Diseases of the skin and subcutaneous tissue
- L: Others

INSIGHT

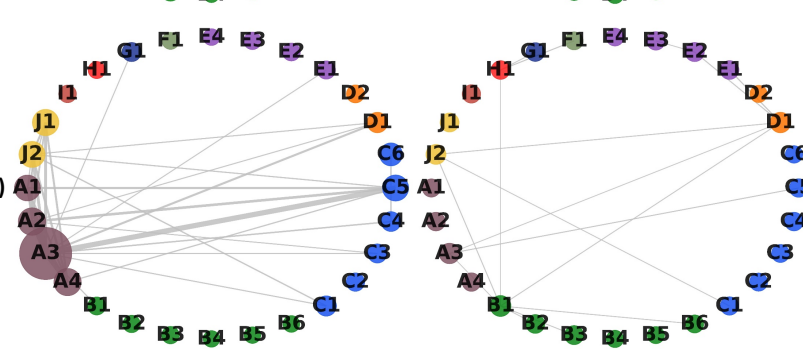
Subphenotype 1
(Cardiac and Renal)
N=7,047



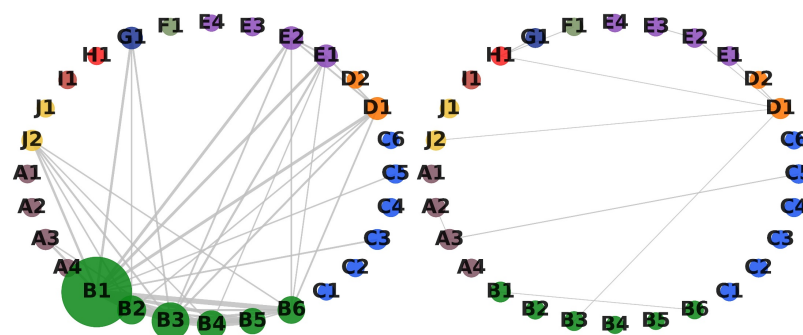
Subphenotype 2
(Respiratory, Sleep and Anxiety)
N=6,838



Subphenotype 3
(Musculoskeletal and Nervous)
N=4,879



Subphenotype 4
(Digestive and Respiratory)
N=2,117



COVID-19 Positive

COVID-19 Negative

A Diseases of the Musculoskeletal System and Connective Tissue

- A1: Osteoarthritis
- A2: Spondylopathies
- A3: Musculoskeletal pain
- A4: Connective tissue disease

B Diseases of the Digestive System

- B1: Abdominal and pelvic pain
- B2: Gastrointestinal disorder
- B3: Esophageal disorder
- B4: Gastritis and duodenitis
- B5: Stomach disorder
- B6: Nausea and vomiting

C Diseases of the Nervous System

- C1: Cognitive problems
- C2: Sleep disorder
- C3: Headache
- C4: Malaise and fatigue
- C5: Nervous system pain
- C6: Nervous system disorders

D Diseases of the Respiratory System

- D1: Breathing abnormalities
- D2: Lower respiratory disease

E Diseases of the Circulatory System

- E1: Chest pain
- E2: Abnormalities of heart beat
- E3: Cardiac dysrhythmias
- E4: Heart failure

F Diseases of the Genitourinary System

- F1: Renal failure

G Diseases of the Blood

- G1: Anemia

H Endocrine, Nutritional and Metabolic Diseases

- H1: Fluid/electrolyte disorders

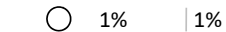
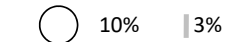
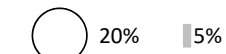
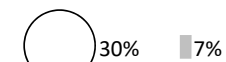
I Mental and Neurodevelopmental Disorders

- I1: Anxiety

J Others

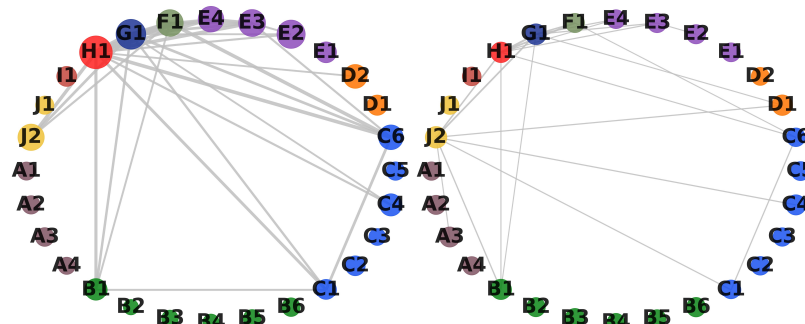
- J1: Skin sensation problems and rash
- J2: General signs and symptoms

Node size Line width

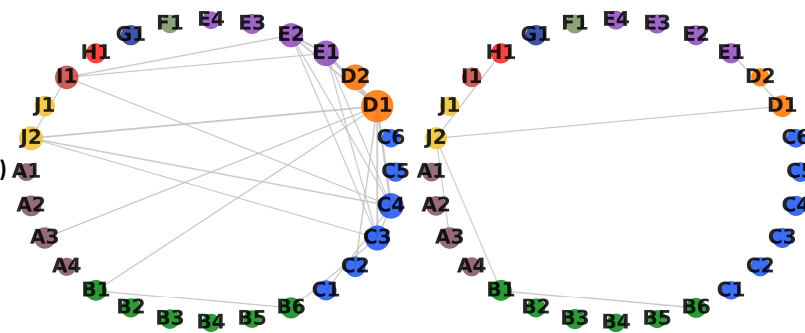


OneFlorida+

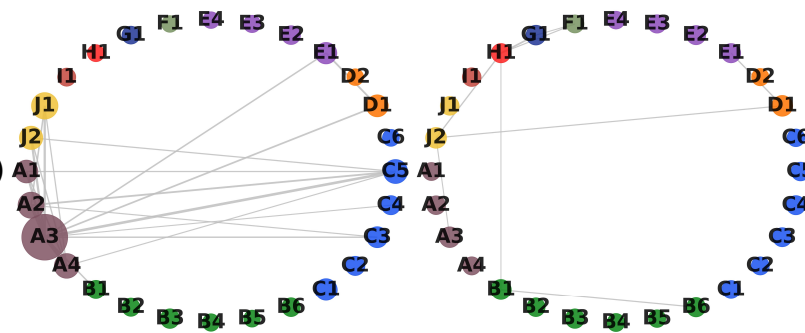
Subphenotype 1
(Cardiac and Renal)
N=3,490



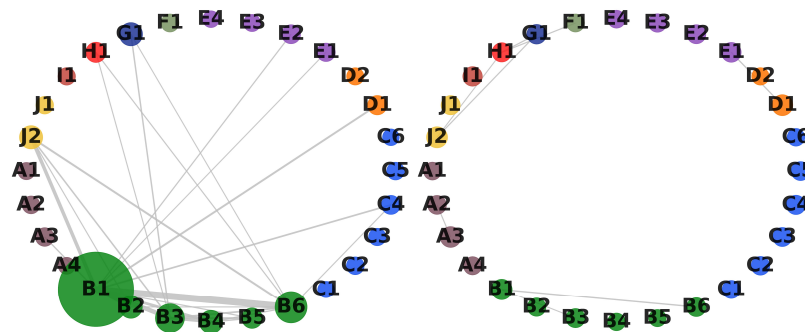
Subphenotype 2
(Respiratory, Sleep and Anxiety)
N=5,281



Subphenotype 3
(Musculoskeletal and Nervous)
N=3,205



Subphenotype 4
(Digestive and Respiratory)
N=1,748



COVID-19 Positive

COVID-19 Negative

A Diseases of the Musculoskeletal System and Connective Tissue

- A1: Osteoarthritis
- A2: Spondylopathies
- A3: Musculoskeletal pain
- A4: Connective tissue disease

B Diseases of the Digestive System

- B1: Abdominal and pelvic pain
- B2: Gastrointestinal disorder
- B3: Esophageal disorder
- B4: Gastritis and duodenitis
- B5: Stomach disorder
- B6: Nausea and vomiting

C Diseases of the Nervous System

- C1: Cognitive problems
- C2: Sleep disorder
- C3: Headache
- C4: Malaise and fatigue
- C5: Nervous system pain
- C6: Nervous system disorders

D Diseases of the Respiratory System

- D1: Breathing abnormalities
- D2: Lower respiratory disease

E Diseases of the Circulatory System

- E1: Chest pain
- E2: Abnormalities of heart beat
- E3: Cardiac dysrhythmias
- E4: Heart failure

F Diseases of the Genitourinary System

- F1: Renal failure

G Diseases of the Blood

- G1: Anemia

H Endocrine, Nutritional and Metabolic Diseases

- H1: Fluid/electrolyte disorders

I Mental and Neurodevelopmental Disorders

- I1: Anxiety

J Others

- J1: Skin sensation problems and rash
- J2: General signs and symptoms

Node size Line width



30%



7%



20%



5%



10%



3%



1%



1%

Outline

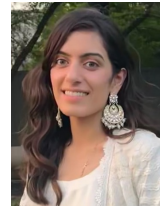
- Introduction
- Subphenotyping of COVID-19 at infection confirmation
- Subphenotyping of Severe COVID-19 after Mechanical Ventilation
- Subphenotyping of Long COVID
- **Discussions**

Discussions

- Complicated diseases are heterogeneous
 - Snapshot
 - Longitudinal
- Identification of disease subphenotypes from patients' clinical data can help us better understand the clinical heterogeneity and trigger stratified medicine
- The next step is to “validate” the subphenotypes
 - External validation on independent data
 - Mechanism investigation
 - Treatment response assessments

Acknowledgements

- NSF
 - 1650723
 - 1716432
 - 1750326
 - 2027970
- NIH
 - R01AG080624
 - R01AG080991
 - R01AG076448
 - R01AG076234
 - RF1AG072449
 - R01MH124740
 - R01MH112148
 - R01GM105688
 - OT2HL161847
- PCORI
- ONR
- NMRC
- Michael J Fox Foundation
- Boehringer Ingelheim
- Google research award
- Amazon
- Sanofi
- MITRE
- Cornell



Lab Manager

- [Zehra Abedi](#)



Instructors

- [Chengxi Zang](#)



Research Associates

- [Zhenxing Xu](#)



Postdoctoral Associates

- [Zilong Bai](#)
- [Daoming Lyu](#)
- [Qianqian Xie](#)
- [Weishen Pan](#)



PhD Students

- Daniel Adler (IS, Together with Deborah Estrin and Tanzeem Choudary, NSF Fellow)
- Matthew Brendel (PBSB, Together with Iman Hajirasouliha)
- Jaqueline Maasch (Computer Science, Together with Volodymyr Kuleshov, NSF Fellow)
- Suraj Rajendran (Tri-I CBM, Together with Iman Hajirasouliha, NSF Fellow)
- Yingheng Wang (Computer Science)
- Manqi Zhou (CB)

<https://wcm-wanglab.github.io/index.html>

Thank You!

few2001@med.cornell.edu

 @feiwang03



<https://wcm-wanglab.github.io/index.html>